

Sincerity and Insincerity

Neri Marsili

1. Introduction

[The study of communication] is in principle the [study of] everything which can be used in order to lie. If something cannot be used to tell a lie, conversely it cannot be used to tell the truth: it cannot in fact be used “to tell” at all.

Umberto Eco, *A theory of Semiotics*

Communication is fundamental to our lives – arguably, it is one of the things that makes us distinctively human. This ability to communicate, however, comes at a price: it makes it possible for speakers to misrepresent what they believe. While communicating allows us to exchange information, it equally makes us vulnerable to misinformation. Communication, then, is a double-edged sword: the ability to communicate sincerely is essentially entangled with the ability to communicate insincerely. The purpose of this book is to explore these two sides of communication: sincerity and insincerity.

I begin by tackling definitional questions. What is sincerity, and what insincerity? How do they differ from related concepts, like truthfulness, deception, omission, or lying? Two main paradigms dominate the theoretical landscape. Speaker-centred views conceive of sincerity as a match (or mismatch) between what the speaker believes and what they claim. Hearer-centred views focus instead on the intended effect on the hearer (informing or misinforming the recipient). Each view admits variations, revisions and expansions. The first half of the book examines how different definitions of sincerity can deal with increasingly complex puzzles, like graded beliefs, non-assertoric speech, and non-literal expressions.

The second half explores normative questions. Philosophers consider sincerity to be an extremely important virtue, and insincerity an unforgivable vice. For some “absolutist” thinkers like Augustine and Kant, no circumstances can ever justify being insincere – not even saving a life. But why do philosophers deem sincerity so important? Answers vary, but they typically invoke our inability to communicate effectively, to maintain meaningful social bonds, and to share knowledge without a

norm of sincerity. I will conclude by analysing sincerity both as a norm and as a virtue, focusing on sincerity's epistemic value, and its role in sustaining our ability to share knowledge through testimony.

2. What is sincerity?

In 1972, operatives connected to Richard Nixon's re-election campaign broke into the Democratic National Committee headquarters at the Watergate complex in Washington D.C., to wiretap phones and photograph campaign documents. Nixon repeatedly denied involvement or knowledge of the break-in, including in a televised address to the nation in 1973, in which he stated:

- (1) I had no prior knowledge of the Watergate break-in; I neither took part in nor knew about any of the subsequent cover-up activities.

Nixon's denials were soon proven untruthful, when Oval Office tapes revealed Nixon's direct involvement in the cover-up. These revelations eventually led to his resignation on August 9, 1974, to avoid impeachment.

Nixon's statement is a paradigmatic case of insincerity: (1) is a believed-false statement, uttered with the intention of deceiving its recipients. Paradigmatic cases of sincerity, by contrast, involve stating what you believe to be true, with the intention of passing information onto your audience. Unlike insincere statements like (1), sincere statements rarely make for great stories: since sincerity is the rule, the fact that someone chose to speak sincerely is typically uneventful.

The adjectives "sincere" and "insincere" can be used in at least two senses. First, it is used to describe a *property of utterances* – that is, of communicative acts made in a certain context. I may say, for example, that (1) was insincere, because in uttering (1) Nixon intentionally misrepresented his involvement in the Watergate break-in. In this sense, insincerity is predicated of a *particular utterance*, relative to a particular agent (the speaker) and a particular context. I shall call this conception that focuses on the sincerity of particular utterances "sincerity in discourse" (or "discursive sincerity"¹), and it will be the main concern of this book.

Second, ascriptions of sincerity are sometimes focused on particular speaker, rather than specific communicative acts. I may say, for example, that Nixon was an insincere politician, and that (1) was definitely in character. When "sincere" and "insincere" are used in this way, they refer to *dispositions* or *character traits* – they

¹ "Discourse" here stands for communication in general. It's less precise, but makes for a slightly more readable label.

describe what a certain person tends to do under certain circumstances (a disposition) or and the psychological features that underlie these dispositions (a character trait). I shall call this the *dispositional* conception of sincerity and insincerity. While I'll occasionally discuss this notion, (§3.3), it will not be the main concern of this essay.

The contrast between dispositional and discursive conceptions of sincerity is lexicalised in our vocabulary for lying. We distinguish between a liar (someone who has a disposition to lie) and a lie (an insincere utterance). Sincerity and insincerity, however, do not admit nominalisation: they cannot be turned into nouns. Still, we can distinguish between an *insincere statement* (close, but not equivalent, to the notion of a lie) and an *insincere speaker* (close, but not equivalent, to the notion of a liar). However, even these compound constructions maintain a level of ambiguity: ascribed to the speaker, insincerity can refer both to a speaker's disposition in general (e.g. Nixon was an insincere politician) or a speaker's communicative act (e.g. Nixon was insincere *in stating* (1)). To disambiguate, I will instead distinguish between *discursive* and *dispositional* conceptions of sincerity, focusing on the former.

2.1 Sincerity in discourse and beyond

2.1.1 Communication: a precondition for sincerity

There is a prejudice against the spoken lie, but none against any other; and by examination and mathematical computation I find that the proportion of the spoken lie to the other varieties is as 1 to 22,894. Therefore the spoken lie is of no consequence, and it is not worth while to go around fussing about it and trying to make believe that it is an important matter

Mark Twain, *My First Lie, and How I Got Out of It*

In his *Parerga e Paralipomena*, Arthur Schopenhauer writes that “there is in the world only one mendacious and hypocritical being, namely man. Every other is true and sincere, in that it frankly and openly declares itself to be what it is and expresses itself as it feels”. Unlike animals, whose ingenuity and transparency are fascinating to us, the degenerate human tendency to lie “stands as a blot on Nature”.

Schopenhauer's harsh comments are comically inaccurate: plants and animals are capable of incredibly creative and complex forms of deception. However, there's a grain of truth in his observation: although many non-human organisms are able to

deceive, there's a sense in which lying and insincerity (unlike deceiving) are exclusive to the domain of coded, linguistic communication.

Take the orchid *Cryptostylis erecta*, which is pollinated by the so-called orchid dupe wasp (*Lissopimpla excelsa*). Male wasps mistake the orchid for female wasps, and copulate with them. While there is a sense in which the wasp is deceived by the orchid's shape, it would be a stretch to say that the orchid *lied* to the wasp, or that the orchid was *insincere*. Instead, it seems appropriate to distinguish between the broader concept of *deception* and the narrower concept of *insincerity*. Anything designed² to induce false representations of the world (even an orchid's appearance) is a form of deception. Insincerity, as I'll understand it, is instead a property of communicative acts. This rules out non-verbal deception, whose best illustration is perhaps offered by Kant's (Kant LE) example of the suitcase packer who ostentatiously prepares their suitcases to persuade a friend that they are about to leave for a long trip.

The idea that insincerity requires linguistic communication can be traced back to Augustine (*Contra Mendacium* XII), who wrote that "a lie is a false signification by words". Aquinas (ST, q110) clarifies that "words" should be understood in a loose sense:

When it is said that "a lie is a false signification by words," the term "words" denotes every kind of sign. Wherefore if a person intended to signify something false by means of signs, he would not be excused from lying.

Indeed, we can communicate insincerely without using words, by using any sort of conventional signals. I might nod my head to express agreement insincerely (i.e. even if I disagree), or I might send insincere messages by using other non-verbal conventional codes, like smoke signals, morse code, maritime flag signals, and so forth.

I shall say, then, that (in our intended sense) sincerity and insincerity are properties of communicative acts. Specifically, my focus will be *intentional linguistic communication*, that is, communication that is achieved by means of conventional signs, and intentionally. I won't attempt to provide a more precise definition of communication here, as it would lead us astray (although I'll draw some boundaries in the next sections). For current purposes, it's enough to have established a first intuitive distinction between non-verbal deception (illustrated by Kant's suitcase

² This raises the question of what exactly qualifies as a deceptive "design". This question is amply discussed in the context of the study animal signalling (see e.g. Smith and Harper 2003, Searcy and Nowicki 2005).

packer and the misleading orchids) and insincere discourse. The next sections will refine this distinction.

2.1.2 Omission

In the movie *The invention of lying*, Ricky Gervais imagines a world in which lying has not been invented yet: everybody is sincere all the time. Not only the characters say only what they believe to be true: they seem unable to withhold information that we would typically keep to ourselves (embarrassing details, sexual desires, etc.). For example, we see a waiter welcoming a couple by disclosing irrelevant personal information: “I’m very embarrassed to work here. And [to the woman] you’re very pretty”. The movie implicitly assumes that omitting information amounts to lying.

While this assumption yields good comedic effects, it also engenders conceptual confusion. I argued that lying rather requires communicating something. By refraining from communicating, then, one cannot lie. Similarly, if we understand insincerity as a property of communicative acts, failing to disclose information doesn’t qualify as insincerity either.

This is not to say that *no* act of omission ever deserves to be called insincere. In some special circumstances, by staying silent a speaker can communicate a defined proposition. For example, an informer wearing a wire might be instructed “Cough if the mafia boss is there, otherwise stay silent”. Here a message is sent by not speaking. Generalising, when there is an explicit mutual understanding that failing to speak amounts to communicating a specific proposition (e.g. “the boss is not here”) we have what I’ll call “silent communication”. *Silent communication* can be insincere, but is clearly different from withholding information: the former involves communicating something, the latter does not (Mahon 2015)³.

That noted, by failing to speak up one can still be culpable of *deception by omission*, which is *pro tanto* morally objectionable. For example, I might resent a friend for keeping me in the dark about my partner’s love affairs, or a colleague for not warning me that they accidentally dropped coffee on my laptop while I was on my lunch pause. In such cases, however, failing to disclose some content *p* doesn’t amount to communicating that *p* is false – at most, it is to let someone *infer* that *p* is false. It would, therefore, not be insincerity *stricto sensu*, since no communicative act is involved.

While deceptive omissions fall out of the scope of “discursive sincerity”, there’s admittedly *some* sense in which they can be insincere. The movie characters in *The Invention of Lying*, who disclose all information that they deem relevant to the conversation, display a propensity to be transparent that closely resembles sincerity,

³ Not all cases are so straightforward; for discussion, Fallis (2018) and §2.4.

despite exceeding its demands. I shall say that they are “supersincere”, in the sense that, on top of saying only what they believe, they also disclose all information that they deem relevant to the audience. Supersincerity sometimes clashes with other social norms (such as politeness and privacy, cf. §3.2), and that’s why Gervais’ movie characters, who can’t help violating such norms, are awkward and comical. Generalising, supersincerity tracks a social expectation that, though akin to sincerity in discourse, is importantly different from it.

2.1.3 Sincerity beyond communication

*I sincerely believe that banking establishments
are more dangerous than standing armies*

Thomas Jefferson

Since this book deliberately focuses on discursive sincerity (and insincerity), it’s worth commenting on the fact that in ordinary language, sincerity is sometimes predicated beyond the realm of communication. Actions can be described as sincere and insincere: we speak of a sincere kiss, a sincere attempt to help. Beliefs, too, admit such descriptions, like Jefferson’s “sincere belief” that banking establishment are more dangerous than standing armies. Even a philosophical doctrine can be described as sincere and insincere. Writing about Schopenhauer’s life, Bertrand Russell (1946) notes:

[His] doctrine [was not] sincere, if we may judge by Schopenhauer’s life. He habitually dined well, at a good restaurant; he had many trivial love-affairs, which were sensual but not passionate; he was exceedingly quarrelsome and unusually avaricious. On one occasion he was annoyed by an elderly seamstress who was talking to a friend outside the door of his apartment. He threw her downstairs, causing her permanent injury

What made Schopenhauer’s doctrine “insincere”? Presumably, a certain inconsistency between what Schopenhauer preached (“the virtue of asceticism and resignation”) and his actions. When Russell wrote that Schopenhauer’s philosophy was insincere, then, he meant that there was a “mismatch” between what Schopenhauer preached and the way he acted. Both in discourse and beyond discourse, at the core of the notion of (in)sincerity there seem to be a contrast

between two states: there's insincerity when there's a "match" between two states, and insincerity when there is a mismatch or disconnect.

In *discourse*, the contrast is between what one says and what one thinks: a sincere statement is one that "matches" one's beliefs. Sincerity in *action* can take a variety of forms. In Schopenhauer's case, the contrast is between what he preaches (his ascetic doctrine) and what he does (his indulging lifestyle). But sincerity in action can also involve a contrast between one's actions and what one *believes* or feels – similarly to sincerity in discourse. A sincere kiss is one where the action (kissing) reflects my inner state (an emotion). A sincere attempt to help is one action that reflects a genuine will to be helpful. Whether the action is sincere depends on whether you have the internal state that is supposed to accompany the action under consideration.

Sincerity in belief is more complex. When Jefferson claims that he "sincerely believes" that banking establishments are dangerous, the expression is idiomatic: Jefferson actually means that he is *sincerely asserting* his claim. But sometimes beliefs are described as "insincere" in a more substantive sense. Paradigmatically, this happens when a believer is somewhat aware that their evidence doesn't support a given belief, but they hold onto their belief nonetheless.

Suppose that Johnny firmly believes that he is a loving husband, despite being aware that he's culpable of all sorts of abuses. It might be said that Johnny's belief that he's a loving husband is "insincere", because he grasps onto this thought despite overwhelming evidence to the contrary. Once again, insincerity stems from a mismatch or a disconnect: Johnny's awareness that his behaviour is abusive, and hence what he should believe, doesn't match what he consciously believes.

Like sincerity in discourse, then, sincerity beyond discourse involves a mismatch or a disconnect, typically between an internal representation of the world (a belief, desire, etc.) and another, external or otherwise (in speech, in belief, or as expressed by an action)⁴. Despite the important analogy between these conceptions, however, sincerity in discourse constitutes a coherent, distinctive phenomenon that deserves treatment in its own right. Now that its cognates (*sincerity in action*, *sincerity in belief*, *supersincerity*) have been briefly introduced and discussed, I shall leave them behind, to focus on sincerity in communication.

2.1.4 Falsity, Truth, and truthfulness

⁴ Arguably, a more appropriate label for sincerity beyond discourse is *authenticity*: many of the examples discussed, such as behaving in ways that don't reflect our true nature and our deepest beliefs, are best described as failures to be authentic, rather than sincere. For an introduction to the philosophy of authenticity, see Varga and Guignon (2023).

It's May 1, 1486, in Córdoba, Spain. Christopher Columbus stands before Queen Isabella I and King Ferdinand II, passionately presenting his ambitious plan. He declares:

- (2) Your Majesties, with a few ships we can reach the Indies by sailing west across the Atlantic Ocean, for about 750 leagues

Famously, Columbus was wrong: he grossly miscalculated the distance between Spain and the Indies. Following his route, it would have taken approximately 3,500 leagues (about 20,000 km) to reach his destination. Columbus' statement was⁵ sincere, but nonetheless false.

The example illustrates a familiar way in which sincerity and truth can come apart: honest mistakes. Our statements can fail to be true even if they are sincere: unfortunately, we humans aren't endowed with the gift of omniscience. Sometimes we get things wrong; when we sincerely assert our wrong opinions, we make statements that are sincere but false.

The opposite can happen, although it's certainly a less familiar case: a statement can turn out to be true even if it's insincere. Consider the following example:

Cheating judo

Giacomo and Frida are married, but Giacomo is looking for an affair. At a party, he starts flirting with a woman called Ana. Ana knows that Giacomo is married. Hoping to make her feel comfortable about flirting with a married man, Giacomo tells her:

- (3) *My wife has been cheating on me in the last months*

As it turns out, unbeknownst to Giacomo, Frida is really cheating on him.

Giacomo is clearly insincere. But his statement turned out to be true: his wife is truly cheating on him. Just like sincerity is no guarantee of truth, insincerity is no guarantee of falsity.

The purpose of these examples is to impress upon the reader the important conceptual distinction between sincerity and truthfulness. Whether (2) and (3) are true or false doesn't depend on the speaker's subjective beliefs, but rather on *objective* features of the world⁶ (whether the Indies are 750 leagues away, whether Frida is cheating on Giacomo). By contrast, whether (2) and (3) are sincere or insincere

⁵ In our imaginary example, that is. To avoid misunderstandings: there is no historical record of Columbus uttering this precise sentence (or its Castilian equivalent).

⁶ All talk of objectivity is rejected by some philosophers. I will leave such controversies aside; see Williams (2002, chap. 1) for discussion.

depends on the speaker's *subjective* mental states (did they believe that what they said is true? Were they aiming to deceive their interlocutor?). To stress the contrast between these two dimensions of evaluation, some philosophers distinguish between truthfulness (objective dimension) and sincerity (subjective dimension), and (correspondingly) untruthfulness and insincerity.

Of course, this terminology is somewhat stipulative. In ordinary language we occasionally use “truthful” and “untruthful” as synonymous with sincere and insincere. Similarly, “telling the truth” is often used as a synonym of “stating one’s opinion”, as opposed to its literal meaning (to state what is objectively true, regardless of whether it matches one’s opinion). In line with these uses, philosophers (e.g. Lewis 1975) sometimes use “truthfulness” in a subjective sense: a statement is truthful if the speaker *aims to tell the truth*, independently of whether they succeed. In this book, I will stipulate that “truthfulness” and “untruthfulness” designate the objective dimension of evaluation of a statement: whether it matches reality. So understood, the opposition between truthfulness and sincerity help us navigate the different ways in which statements can be epistemically defective.

Strong of these conceptual distinctions, let’s move on to analyse two main conceptions of sincerity: one that defines it as a relation between the speaker’s statement and their mental states (speaker-centred conception, §2.2), and one as a relation between their mental states and the mental states that they aim to induce in the hearer (hearer-centred conception, section §2.3).

2.2 The expression view: speaker-centred sincerity

*Hateful in my eyes, even as the gates of Hades,
is that man that hideth one thing in his chest and sayeth another
(Achilles, in the Iliad, IX 312-31)*

According to an influential philosophical view, communication is essentially a tool for self-expression – that is, for making public what’s inside our minds. Without speech, every individual is an island. Our internal attitudes (beliefs, intentions, emotions) are private states inaccessible to others. The power of communication is that it allows us to make these states public: to *express* our mental states.

The idea that words are the outer expression of mental states can be traced back to Aristotle (*De Interpretatione*, § 1). It was later endorsed by Locke (1690, III, II, §1), and then revived and systematised in contemporary philosophy of language (e.g. Frege 1956; Searle 1979; Davis 2003; Green 2007a). Inevitably, the notion of attitude expression is understood in slightly different ways by these authors. Broadly put, for a speaker to *express* a psychological state is for the speaker to *represent themselves* as being in that psychological state. Put more bluntly: you express a mental state when

you communicate that you have that mental state. For our purposes, I don't need to settle on a precise theory of thought-expression – that's for another book. But it's important to stress what expression, as it's understood here, *does not* involve.

First, expression is independent from persuasion (even *attempted* persuasion)⁷. When I say “Ouch!”, my cry expresses pain regardless of whether I am attempting to convince you that I am in pain (my cry might be a mere reflex), and regardless of whether I succeed in convincing you that I am in pain. Second, *expressing* an attitude doesn't entail having that attitude⁸: my shouting “ouch!” can express pain regardless of whether I am, in fact, experiencing pain. Third, saying something is not enough for expressing a belief; the speaker must actually *communicate* the expressed proposition. For example, if I ironically remark “Oh, sure, I'm a Swiftie, big time!” to point out that I obviously don't like Taylor Swift's music, I've *said* that I'm a Swiftie, but I've not expressed the belief that I'm a Swiftie, since there's mutual understanding that I don't mean what I say⁹.

On this conception of communication, which I shall call the *expression view*, sincerity is deeply tied to thought expression. Expressing our private mental states, opens the possibility of misrepresentation. When there is a *match* between what the speaker says and what they believe, the statement is sincere. When there is a mismatch instead, the statement is insincere (Figure 1)¹⁰. The result is the following account of sincerity:

STANDARD VIEW

Sincerity depends on whether the speaker's (S) expressed belief matches their actual beliefs

- *A statement is sincere iff S believes its expressed content p to be true*
- *A statement is insincere iff S believes its expressed content p to be false*

This “classic” or “standard” conception of sincerity as a *match* between speech and thought (and insincerity as *mismatch*) finds resonance in a long philosophical tradition. In the IV century, Augustine defined lying in terms of a *duplicity of the*

⁷ This in contrast with neo-Gricean accounts of expression, like Bach and Harnish's (1979), which require attempted persuasion as a necessary condition for expression.

⁸ This is in contrast to authors who define expression in a way that requires having the corresponding mental state, like Davis (2003, 25), Green (2007a, 70–83), Owens (2006).

⁹ If something, I expressed the belief that I'm not a Swiftie. I'll come back to non-literal communication in §2.4.

¹⁰ Here I only discuss *assertoric* communication, whose sincerity depends on the speaker's *beliefs*. I cover other mental states (intentions, desires, etc.) and other communicative acts (promises, requests, etc.) in §2.7.

*heart*¹¹ (“*duplex cor*”, *De Mendacio*, III, 3): “[the liar] has one thought hidden in his breast and another ready on his tongue, and this is the evil proper to the liar” (*Contra Mendacio*, c. 421). The idea that lying consists in speaking in opposition to one’s mind becomes then standard in medieval philosophy; following Raymond de Peñafort, Clem (2023) calls this the “*contra mentem* principle”.

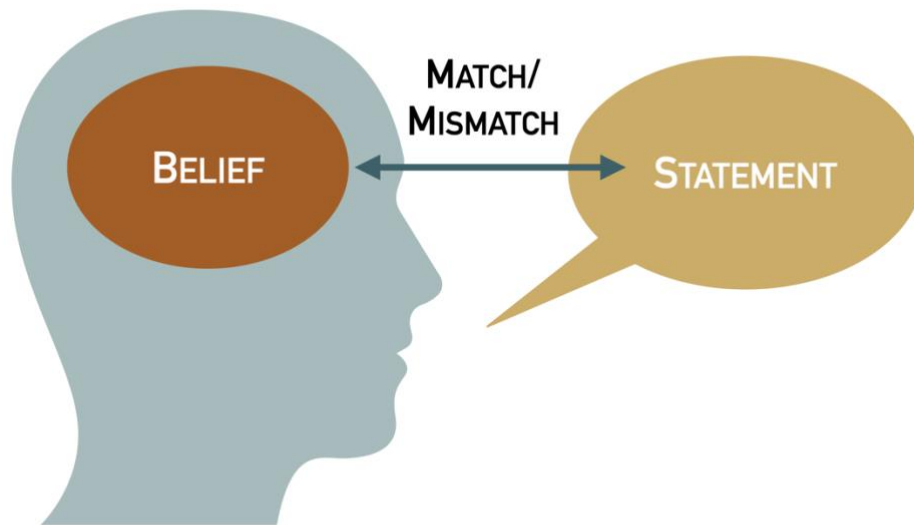


Figure 1: A visual representation of the standard expression view. Sincerity depends on whether the speaker believes or disbelieves the proposition expressed by the statement.

2.2.1 Sincerely asserting what you believe to be false

While STANDARD VIEW works well for paradigmatic cases of sincerity and insincerity, it struggles to accommodate cases involving self-deception and fragmented beliefs¹². To illustrate, consider the following scenario (from Ridge 2006, edited)

LOVELY MOTHER

¹¹ For elaboration of this view, see Griffiths (2004).

¹² There’s substantial scholarly debate on what self-deception is and how it works. For discussion, see Krstić (forthcoming).

Beppe believes that he believes his mother loves him, but deep down he actually doesn't believe that she does. In fact, Beppe believes his mother hates him. Beppe's behaviour betrays that he deludes himself about his own beliefs. For example, he might predict that his mother will do things that would make sense only if he believed she hated him. Beppe feels a sort of anxiety people associate with being hated by a close family member, goes out of his way to try to please her, is highly deferential to her, etc. Nonetheless, because he cannot cope with the idea that his mother hates him, he has somehow convinced himself that he believes that she loves him. When we ask Beppe whether his mother loves him, he replies:

(4) *Yes, of course she does*

Beppe's assertion is intuitively sincere. After all, *he believes that he believes* (i.e. he possesses a "higher order belief" – a belief about what he believes) that his mother loves him. However, what Beppe *really believes* (his first-order belief, independent of what he believes that he believes) is that his mother hates him. Interpreted literally, the standard view classifies (4) as *insincere*. Many scholars take this to be a problem (Stokke 2018, 173–75; Chan and Kahane 2011; Ridge 2006, 488–89).

As a solution, some suggest that the sincerity of a statement depends on the conscious, higher-order mental states of the speaker, rather than their deep, actual beliefs¹³ (see e.g. Mellor 1977; Moran 2005). On this view, the sincerity of a statement determined by the relationship between what one says and their *higher-order beliefs* (their beliefs about what they believe):

HIGHER-ORDER EXPRESSION

Sincerity depends on whether the speaker's (S) expressed belief matches what they believe that they believe

- *A statement is sincere iff S believes that S believes it to be true*
- *A statement is insincere iff S believes that S believes it to be false*

HIGHER-ORDER EXPRESSION gives the right verdict for cases of fragmented belief like LOVELY MOTHER. By its light, Beppe is sincere in asserting (4) even though Beppe actually believes that his mother hates him, because what matters is that he *believes that he believes* that she loves him. Standard insincere statements like (1), too, are classified as insincere, since on top of believing that (1) is false, Nixon

¹³ Higher-order beliefs and conscious beliefs can come apart under some circumstances. I'll ignore this complication here, as well as formulations in terms of *assent* (rather than belief), but see Stokke (Stokke 2018, 177–80) for discussion.

believes that he believes it. Nonetheless, even the HIGHER-ORDER EXPRESSION faces some challenges.

2.2.2 Misspeaking and intentions

How many animals of each kind did Moses take on the ark? If you answered “Two”, you are in good company (Erickson and Mattson 1981). Yet, this isn’t what you actually believe. You know that, according to the Biblical myth, it was Noah, not Moses, who took animals on the ark. This common slip of tongue, known as the ‘Moses Illusion’, illustrates the concept of *misspeaking*: we sometimes try but fail to say what we want to say.

Those who misspeak are clearly not liars (Sorensen 2011). Sure, they fail to say what they believe. But accusations of mendacity aren’t in order here. The victims of the Moses illusion might be distracted, but they are certainly not insincere. This yields a new problem for speaker-centred views. If you accidentally respond “two”, you say something that you neither believe, nor believe that you believe. Every definition reviewed so far misclassifies misspeaking as insincere speech.

If this verdict is intuitively incorrect, it’s arguably because misspeakers have no *intention* to misrepresent their beliefs. Perhaps, then, insincerity is a matter of communicating a proposition that *intentionally* misrepresents one’s belief – and conversely, sincerity a matter of *intentionally* representing one’s beliefs¹⁴:

INTENTIONAL EXPRESSION:

- A statement is sincere iff in making their statement, the speaker (S) intends to express a proposition that S believes to be true
- A statement is insincere iff in making their statement, S intends to express a proposition that S believes be false

There are notable precursors to this view. Augustine (DM) defined lying as an “utterance accompanied by the intention to utter a falsehood”. Similarly, Aquinas (ST, II, II, 110) argues that “the definition of lying is taken from formal falsity, i.e. from the fact that someone has the intention to state what is false”¹⁵. Many

¹⁴ Of course, what matters for sincerity is the speaker’s *conscious* intention. As I understand it, conscious intention is about having reflective *access* to the goals guiding one’s action, rather than explicitly *thinking* about such goals. So, if I instinctively and spontaneously tell the truth, I can be sincere even if I have not explicitly thought “I will now tell the truth”. For further discussion of conscious intentions and sincerity, see Stokke (2018, 181–86).

¹⁵ These formulations require an intention to utter an *objective* falsehood, rather than a believed-false proposition. Still, any rational intention to utter a falsehood entails an intention to utter a believed-false proposition, hence the analogy with the intentional expression view.

contemporary authors, too, have linked sincerity to an intention to express a proposition that matches one's beliefs (or mismatches, for insincerity)¹⁶.

Since victims of Moses illusions (and many other misspeakers alike) have no intention to communicate what they believe to be false, INTENTIONAL EXPRESSION handles these utterances correctly. Interestingly, even though it doesn't mention higher-order beliefs, it also gives the correct verdict for cases of fragmented belief like (4). A speaker who intends to express a proposition they believe to be true will assert their conscious, higher-order, belief, even if this doesn't correspond to their subconscious, or lower-order, attitudes. Consequently, as things stand, INTENTIONAL EXPRESSION offers the best version of the expression view on the market¹⁷.

2.3 The manipulation view: hearer-centred sincerity

According to its detractors, the expression view overlooks importance of the receiving end of the communicative exchange: the audience (aka the hearer, the audience, the recipient). If people communicate, it's often to get other people to think or do something. Perhaps, then, the essence of communication (its essential and primary function, what makes it important and useful) rather resides in its ability to *influence* other people's thinking and acting. I shall call this the *contagion* or *manipulation* model, because it focuses on how mental states can hop from one mind to another (hence contagion), and how communicators can influence each other's thoughts and action, often to their advantage (hence manipulation). If the expression view is *speaker-centred* (communication is about expressing the internal states of *the speaker*), the contagion view is *hearer-centred* (it's about affecting the internal states of *the hearer*).

In animal communication studies, communication is almost invariably conceived as a form of manipulation – as a system of “signal-response” pairs. Take the bright coloration of noxious frogs, which signals toxicity. Its function is to trigger a certain behaviour in the receiver – specifically, to get predators not to eat the frog. Receivers evolved to comply because have something to gain from executing the appropriate

¹⁶ For Chan and Kahane (2011, 229), the speaker's dominant *motivation* determines sincerity. Fallis's (2012, 578) definition of lying closely aligns with Augustine and Aquinas. In Marsili (2017, 29–30; 2018a, n. 8; 2021a, 504; 2021b, 3258) I defend a version of the intentional expression view. Stokke's (2018, 192) view is much more complex, but equally relies on intentions.

¹⁷ The intentional expression view also accommodates some tricky cases discussed by Pepp (2018) – at least, this is what I argue in Marsili (2021a). Stokke (2018) offers a more sophisticated version of this view; on top of an important trade-off in complexity, I worry that it makes some incorrect predictions, as discussed in footnotes XXX and XXX.

response – they avoid poisoning or death. Generalising, the function of animal signals is typically to manipulate the receiver’s behaviour¹⁸.

Ruth Millikan famously extended this model to human communication. In Millikan’s framework, the evolutionary function of assertion (the purpose for which it was selected and passed on from generation to generation – its “proper function”) is to *persuade* the audience that its content is true (Millikan 2005, chap. 8). The function of assertion, on this view, is to manipulate the thoughts of the audience.

Another influential defender of the manipulation model is H.P. Grice. For Grice (1957) and his followers, genuine communication (“non-natural meaning”) arises when agents produce utterances that are intended to alter the mental states of the receiver in certain ways. Specifically, to mean that p is to have a reflexive intention (R-intention) to get the hearer to believe p (at least partly) because of their recognition that we intend them to believe that p ¹⁹. Here, too, communication is understood as an attempt to influence the mental states of the receiver.

Under this family of views, sincere and insincere speech can be distinguished on the basis of their differential (intended) effects on the audience. Sincere speech is about *coordination* or *contagion*: the beliefs of speaker and hearer are meant to match. Insincere speech, by contrast, is about *manipulating* thought (mismatch is the goal) – see Figure 2.

¹⁸ Many complications emerge, especially in modelling deceptive signals (such as those sent by non-noxious frogs with bright coloration). For an introduction, Smith and Harper (2003), Searcy and Nowicki (2005), and Graham (2020).

¹⁹ Grice’s original formulation takes a speaker S to mean p by uttering U iff:

- S uttered U intending that (1) [A believes P]
- S uttered U intending that (2) [A recognises intention]
- S uttered U intending that [(1) is caused (at least partly) by (2)]

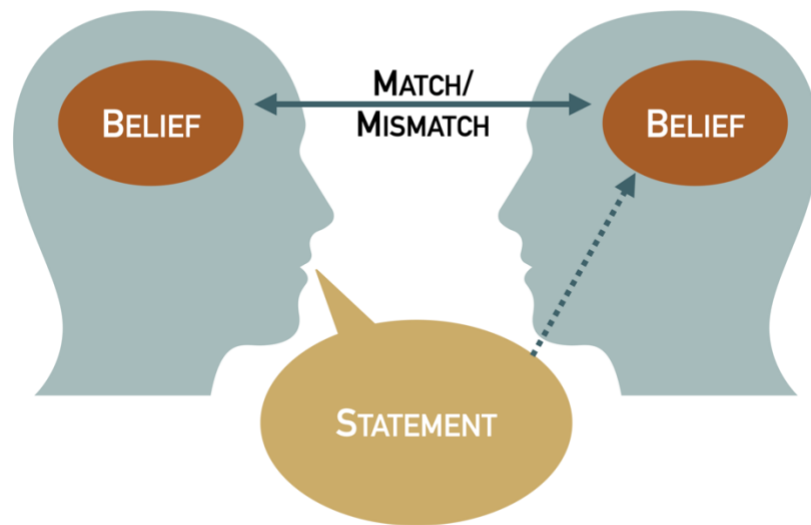


Figure 2: A representation of hearer-centred views, highlighting its focus on the intended outcome: sincerity depends on whether the statement is expected to cause the hearer's belief to match the speaker's.

There are different ways to articulate this idea into a definition; much depends on which precise intention is deemed essential. For instance, the early Grice (1957) understands communication as an attempt to get the audience to believe *that what the speaker said is true*. Subsequent formulations, by both Grice and his followers, rather understood communication as an attempt to get the audience to²⁰ believe *that the speaker believes that what they said is true*. Revisions then escalated in increasingly complex formulations (e.g. Schiffer 1972; Bach and Harnish 1979). Limiting attention to the first two models, we get the following account of sincerity (with its “higher-order” variation in square brackets):

STANDARD HEARER-CENTRED

Sincerity depends on whether the speaker S intends their statement to cause their audience A 's belief to match theirs

- *A statement with content p is sincere iff it's intended by S to cause A to believe [that S believes] p , and S believes p to be true*
- *A statement is insincere iff it's intended by S to cause A to believe [that S believes] p , and S believes p to be false*

²⁰ Or to give the audience a reason to believe that they believe it. More variations are discussed in the literature than I can cover here.

This conception retains important similarities with speaker-centred views: ultimately, whether the speaker believes the relevant proposition determines its sincerity. A speaker that satisfies hearer-centred conditions for sincerity (or insincerity) thereby satisfies speaker-centred conditions, although the opposite isn't true: satisfying speaker-centred conditions doesn't entail satisfying hearer-centred conditions. This view is thus *stronger* than its speaker-centred counterpart – it imposes higher standards for both sincerity and insincerity. Only statements that are intended to cause a certain doxastic²¹ effect on the audience can be deemed sincere or insincere. Contagion theorists will see no problem here: on their view, (assertoric) communication requires such intentions. However, in the next section I'll argue that speakers can be sincere and insincere even when such intentions are absent.

Intermediate positions are possible. Williams (2002) rejects the idea that asserting (sincerely or otherwise) requires an intent to affect the beliefs of the audience, but shares the intuition that *insincerity* requires an intent to modify the audience's beliefs: "sincere assertions do not necessarily have the aim of informing the hearer; but insincere assertions do have the aim of misinforming the hearer" (Williams 2002, chap. 4.2). The result is that one is insincere whenever one is attempting to get the audience to believe some²² false proposition, and sincere otherwise.

DECEPTIVE HEARER-CENTRED

- *A statement is insincere iff it's intended by S to cause A to believe a false proposition*
- *A statement is sincere otherwise*

Introducing some terminology is helpful to appreciate systematic differences between definition. STANDARD HEARER-CENTRED is a "bound" conception of sincerity, because it ties sincerity to a particular proposition: whether the speaker is sincere depends on whether the speaker is attempting to get the audience to believe *a particular proposition*. By contrast, DECEPTIVE HEARER-CENTRED is "unbound", because it doesn't specify which false proposition one has to aim to instil in the audience. Additionally, since sincerity is defined as the absence of insincerity, I'll say that DECEPTIVE HEARER-CENTRED is a "mirrored" definition.

Of course, more definitions can be developed by combining the elements reviewed so far. STANDARD HEARER-CENTRED can be made into a mirrored definition by defining sincerity as the absence of insincerity, or into a bound

²¹ Doxastic is a technical term, meaning "pertaining to belief".

²² Williams rejects the idea that insincerity should be tied to a specific proposition: deceiving about *any* proposition is enough to be insincere.

definition by removing its ties to a specific proposition. The exact intention required, the target proposition, etc. can also be changed²³. The options are so many that attempting to list them all would be hopeless; for our purposes, all that matters is that these definitions admit systematic variations.

Compared to speaker-centred conceptions of sincerity, which are orthodoxy in pragmatics and philosophy of language, hearer-centred conceptions are somewhat fringe views. Still, they have a long history. The earliest discussion of which I'm aware is Aquinas. The *Doctor Angelicus* distinguishes three ways in which assertions can fail to be true²⁴.

[there can be falsehood] materially, since what is said is false, formally, on account of the will to tell an untruth, and effectively, on account of the will to impart a falsehood. (ST. II, 110,1)

Aquinas here differentiates between three notions. The first, material falsehood, is simply *falsity* (or objective *untruthfulness*, cf. §2.1.4) rather than *insincerity*. The second, the formal falsehood, corresponds to the “intentional expression view” discussed in XX. Finally, the last notion is a hearer-centred conception of insincerity: the will to impart a falsehood comes very close to DECEPTIVE HEARER-CENTRED²⁵.

One might worry that hearer-centred views conflate the notion of an *insincere* statement with the notion of a statement that is *intended to deceive*. DECEPTIVE HEARER-CENTRED certainly invites this equivalence²⁶. To intend to deceive is to intend to instil a false belief, which is precisely how this view defines insincerity. If hearer-centred views conflate two notions that speaker-centred views hold separate, there is a worry that our conceptual repertoire will be impoverished by adopting them. Why not use, instead, the term (in)sincerity for speaker-centred sincerity, and intended deception (and lack thereof) for hearer-centred (in)sincerity?

While these considerations are valid, it's also true that we do expect speakers to refrain from deceiving, and that this expectation can reasonably be regarded as

²³ For example, see Eriksson (2011, 230) for a complex variation of STANDARD HEARER-CENTRED.

²⁴ Latin lacks a dedicated word for sincerity, so that Aquinas' discussion of truthfulness (*veracitas*) covers both sincerity and truthfulness at once.

²⁵ Hearer-centred views are also discussed in Trilling (2009, 58), who argues that this distinction is paralleled by the French and English norms of sincerity: simply put, Frenchs demand speaker-centred sincerity; whereas Brits hearer-centred sincerity. While the suggestion is fascinating, the fact that most English philosophers and linguist defend a speaker-centred conceptions of sincerity should give us pause.

²⁶ While STANDARD HEARER-CENTRED does not treat them as fully equivalent, it still allows for a significant overlap between the two notions.

related to sincerity. To adjudicate whether hearer-centred conceptions really track our intuitions, let's bring more considerations to the table, and move on to assess this view's ability to handle specific cases.

2.3.1 Bluffs

*I never could tell a lie that anyone would doubt,
nor a truth that anybody would believe.
Mark Twain, Following the Equator*

Speaker-centred and hearer-centred views have important similarities, and their verdicts converge in most circumstances. But there is space for divergences, such as the following example, which goes back to St. Augustine (here reinterpreted):

TWO ROADS

Simplicius needs to choose between two paths to reach his destination. Marcellus, his business rival, knows that one road is teeming with bandits, while the other one is safe. He also knows that Simplicius will believe the exact opposite of whatever he claims. If Marcellus claims that the road on the right is safe, Marcellus will take the one on the left, and vice versa.

What should Marcellus do? An option is to bluff: Marcellus can claim that the road with the bandits is safe, knowing that Simplicius will believe the opposite, and take the safe road. Here Marcellus will get Simplicius's beliefs to match his, at the price of saying what he believes to be false. To simplify discussion, let's call this option ALTRUISTIC UNTRUTH. If instead Marcellus states what he believes to be true (namely, that the safe road is the safe one), his interlocutor's beliefs won't match his: Simplicius will end up believing the opposite, and be robbed by bandits. I'll call this option SELFISH TRUTH, assuming that Simplicius has some interest in getting his competitor out of business.

For speaker-centred views of all flavours, the ALTRUISTIC UNTRUTH is insincere (since it involves speaking *contra mentem*) and the SELFISH TRUTH sincere. STANDARD HEARER-CENTRED yields the opposite verdict: the ALTRUISTIC UNTRUTH is sincere (since it achieves a doxastic match between hearer and speaker), whereas the SELFISH TRUTH is insincere²⁷. DECEPTIVE HEARER-CENTRED, instead, deems both utterances insincere: in both cases, Simplicius intentionally gets Marcellus to believe

²⁷ Not so for the "higher order" version of STANDARD HEARER-CENTRED, whose predictions align with speaker-centred views instead (if the qualification mentioned in the next footnote holds).

a false proposition: in SELFISH TRUTH, that he believes that what he says is false²⁸; in the ALTRUISTIC UNTRUTH, that what he says is false.

Which characterisation is correct? The ALTRUISTIC UNTRUTH ensures that Simplicius acquires a true belief: accusations of insincerity might be too harsh here. However, in this case Marcellus also expresses a belief that he rejects (that the road is safe): it would seem too generous to deem his statement sincere. Augustine (who was rather concerned with whether the speaker would be *lying*), did not find a definitive solution. Personally, I have the intuition that Marcellus ALTRUISTIC UNTRUTH would be *insincere* after all, and the SELFISH TRUTH *sincere*. Most authors agree (e.g. Moran 2005, 17; Mahon 2015, sec. 1.2; Fallis 2010, 11), but there are exceptions (Chisholm and Feehan 1977, 153–54; Faulkner 2013). A more decisive case against hearer-centred views emerges from blatant lies.

2.3.2 Blatant lies

*E sappi che tu troverai di molti che
mentono, a niun cattivo fine tirando né di
proprio loro utile, né di danno o di
vergogna altrui, ma perciò che la bugia
per sé piace loro, come chi bee non per
sete, ma per gola del vino.*

Giovanni della Casa, *Galateo, ovvero De' Costumi*, XXIII

A blatant (or bald-faced) lie is a lie that cannot be meant to deceive, because there is mutual awareness that the speaker is lying, so that the speaker cannot believe that the audience can be persuaded. Here's an example:

CCTV CAMERA²⁹

Pete took part in a robbery. He knows that his involvement in the crime was unmistakably recorded on CCTV camera, so that there is no chance that anybody will believe him if he denies that he was there. However, he also knows that if he denies being involved in the robbery, the judge will delay the trial and set a low bail. So he claims that he wasn't present at the scene of the crime. He has no intention to deceive anyone. He only says this because he is planning skip bail.

²⁸ Assuming that this is Marcellus' intention.

²⁹ This example, taken from Marsili (2023), is inspired on the "Witness on CCTV" example (Carson 2006, 289–90) and the "Cheating Student" example (Carson 2006, 290).

Intuitively, Pete is lying, even if he lacks an intention to deceive his audience. Pete’s goal in making a false assertion isn’t to persuade the judge, but rather to motivate the audience to act in certain ways. This is but one example of a lie that is not meant to deceive. Countless more have been discussed in the literature – if you don’t find CCTV CAMERA convincing, there’s a strong chance that you will find that some other demonstrates the point³⁰. The growing list of counterexamples led most scholars to conclude that lying typically, but not necessarily, involves deception³¹.

Related observations have been brought up against manipulation views in general. Gricean accounts have taken quite a beating in the literature, and have somewhat fallen out of fashion in speech act theory, primarily because of their inability to account for “non-manipulative” assertives – that is, claims that are not aimed at persuading the audience (e.g. Aldrich 1966; Alston 2000, 44–50; Siebel 2003; 2020; García-Carpintero 2004; Green 2007b, 75–82). Amendments have been proposed³², but these patches only manage to cover a portion of the objections raised. Wayne Davis (1999) offers a pretty damning overview of the counterexamples faced by this view:

For example, if [the speaker] were proclaiming his innocence in the face of a mountain of incriminating evidence, affirming his beliefs before an inquisitor is forcing him to recant, talking to someone he knew did not trust him, uttering a platitude, answering a rhetorical question, reminding someone of an appointment, answering a teacher. The speaker may not even expect his audience to understand him. He may wantonly baffle his audience by using obscure vague ambiguous, esoteric or foreign language. Alternatively, the audience may be unperceptive, unintelligent, unconscious, or even dead, as when people speak to babies, pets, and the recently departed. More radically, there may not even be an intended audience, as when one is recording something in a private diary, scribbling notes to solve a problem, or venting frustrations by cursing loudly precisely because no one can hear

These are assertions that can clearly be sincere or insincere depending on context. However, for STANDARD HEARER-CENTRED they can be neither (since they cannot be intended to persuade the audience), and for DECEPTIVE HEARER-

³⁰ Among them, lies under coercion (Siegler 1966, 129), ‘bald-faced lies’ (Carson, Wokutch, and Murrmann 1982; Sorensen 2007), ‘knowledge lies’ (Sorensen 2010), ‘tell-tale sign lies’ (Krstić 2019), ‘alternative motivation’ lies (Rutschmann and Wiegmann 2017; Sneddon 2020), and many others (Marsili 2016, sec. 9.2; 2021b, sec. 2.3; Sorensen 2018)

³¹ Some have resisted this conclusion, typically criticising the validity of some examples. For an overview, see Mahon (2015) and Krstić (2023).

³² See, for example, Grice (1989, chap. 5), Strawson (1964), Schiffer (1972), Bach and Harnish (1979).

CENTRED they cannot be insincere (since they cannot be intended to persuade the audience of something false).

Recapitulating, contagion views face two related problems. On the sincerity side, speakers who assert do not always aim to persuade their audiences: STANDARD HEARER-CENTRED incorrectly predicts that such speakers are neither asserting nor being sincere. On the insincerity side, liars don't always intend to deceive, but all hearer-centred views predict that such non-deceptive lies are sincere. Crucially, no view that avoids both objections would be truly hearer-centred – since it would not truly require that the hearer's mental states are affected by the statement.

This is not to say that the contagion view is completely unhelpful for theorising about insincerity. There are, on the contrary, several lessons to take home. Even if assertoric communication is even if not *exclusively* aimed at belief-contagion, this is its *primary* or *typical* goal. Correspondingly, the contagion view identifies *paradigmatic* cases of sincerity and insincerity: the typical liar is one that aims to deceive, and the typical sincere assertor one that aims to inform their interlocutor. There's a corollary: where the predictions of speaker-centred and hearer-centred views diverge, we have non-paradigmatic cases of sincerity and insincerity. A sincere speaker that doesn't aim to persuade their audience is still making a sincere assertion, but their intentions are somewhat atypical; same for the liar who lacks an intention to deceive.

Additionally, it seems that a liar who intends to deceive is *more insincere* than a liar who hopes they won't be believed, or that is indifferent about the outcome. If this is right, hearer-centred views capture something about insincerity that speaker-centred views fail to grasp. Even though deceptive intent isn't necessary for insincerity, deceptive intentions affect the *degree* or *extent* to which the statement is insincere³³. Far from simply being on the wrong track, then, hearer-centred conceptions complement our understanding of what sincerity is. Appropriately incorporated (as an account of paradigmatic cases), they help us appreciate how insincere a speaker can be, introducing more nuance to the picture delineated by the expression view.

Both speaker-centred and hearer-centred view admit further refinement, to accommodate phenomena like non-literal meaning, non-assertoric speech, and uncertainty. The next sections will deal with such complications.

2.4 Literally sincere, indirectly deceptive

³³ I doubt that the same can be said of sincere speech: intending to persuade the audience does not make your sincere assertion more sincere.

In a landmark case for perjury statute, *Bronston v. United States*, (409 US 352, 1973), the defendant (Bronston) was asked if he ever owned a Swiss bank account. He responded:

(5) The company had an account there for about six months, in Zurich.

In fact, Bronston himself owned a bank account in Switzerland. Although (5) is literally truthful, it's also meant to convey something false, namely that Bronston (rather than the company) never had a Swiss account. The US Supreme Court eventually ruled that Bronston's statement, being literally true, did not amount to perjury. However, it would be a stretch to say that Bronston answered the question sincerely. If this intuition is on the right track, insincere speech extends beyond literal communication.

In the specialised literature, non-literal statements that indirectly convey a believed false-proposition are said to be misleading, but not lies (Saul 2012)³⁴. The distinction between lying and “merely misleading”, in turn, is typically understood to be grounded on the distinction between *what is said* and *what is implicated*³⁵ (or the related distinction between *literal* and *non-literal* content, *explicit* vs *implicit*, *asserted* vs *implied*, etc.). Apart from theoretical interest, the distinction has practical implications: as the example illustrates, it's important for legislative purposes (S. P. Green 2018), and is said to have moral (Adler 2006; Strudler 2009; Saul 2012) and epistemic (A. Green 2017; González De Prado 2023) implications.

An advantage of indirect deception (misleading) over lying is that it strategically preserves plausible deniability. A speaker who *merely implies* (without explicitly asserting) a deceptive message is able to deny without contradiction that they meant to convey the deceptive content. If accused by the prosecution of having claimed that he *personally* owned a Swiss account, Bronston could've insisted “that's not what I meant, you misunderstood me: I said that *the company* had an account there”. Deniability is sought by deceivers because it renders it more difficult to hold them accountable for what they communicated (Pinker, Nowak, and Lee 2008; Mazzarella 2023; Saul 2024).

As anticipated earlier, accounts of sincerity can be distinguished based on whether they are “bound” and “unbound”. *Bound* views tie the sincerity to the explicit, semantic content of the statement. On these views, merely misleading

³⁴ This distinction is standard, but not without its detractors. Some reject it altogether (Meibauer 2005; 2014), others argue that it does not fully overlap with the distinction between saying and implicating (Viebahn 2017; 2021; but cf. Marsili and Löhr 2022).

³⁵ The noun “implicature”, and the adjective “implicated” are technical terms introduced by Grice (1989). For a primer, Davis (2010).

utterances like (5) are sincere. Since Bronston believes that the literal content of his assertion, he is sincere by *bound* speaker-centred standards. *Unbound* views, by contrast, don't tie sincerity to a specific propositional content, meaning that indirect statements can be classified as insincere.

Should sincerity be modelled along bound or unbound lines? Calling misleading statements like (5) “sincere” may sound hopelessly lenient. Academics haven't much considered the issue, and only a few scholars explicitly acknowledge that misleading statements fall under the remit of insincere speech (Eriksson 2011, 230; Stokke 2018, 190–93). If one feels this way, any bound view can be converted into its unbound counterpart, with some simple adjustments. If a view of sincerity requires that the literal content p meets condition C, its bound version can be derived by requiring instead that *all the propositions the speaker intentionally communicates* (both the literal and the non-literal ones) meet condition C (otherwise, the statement is insincere). Applied to standard expression view, for example, we obtain the following unbound (mirrored) view:

UNBOUND STANDARD EXPRESSION

Given a speaker S that, in making a statement, communicates³⁶ the set of propositions P

- *The statement is sincere iff S believes all members of P*
- *The statement is insincere otherwise*

This gets the desired results. Misleading statements like (5), as well as canonical lies, are classified as insincere, since both communicate believed-false propositions. To be sincere, the speaker has to believe both the literal and non-literal contents to be true. The same formula can be applied to convert any bound account of insincerity into an unbound one³⁷.

Bound sincerity is always “stricter” than its unbound counterpart, since all the communicated propositions (not just the literal content) must meet the sincerity condition (whatever we take it to be). An advantage of unbound view is that they seem to better track ordinary language use: intuitively, Bronston's statement isn't fully sincere. Additionally, these views display sensitivity to another way in which insincerity can be incremental. Arguably, the more insincere propositions the speaker communicates, the more insincere their statement is. Just like the presence

³⁶ For full symmetry with expression views, one would have to require that the speaker *expresses* a belief in each of these propositions. I chose this formulation to avoid controversy on whether the notion of “expression” applies to indirect speech. What is meant by “communicate” is itself a tricky question, but one too wide to be tackled in this short book.

³⁷ Similarly, by the same token, any unbound view can be converted into a bound one, by tying sincerity *only* to the literal content of the utterance.

(or absence) of a deceptive intent can increase (or decrease) how insincere the statement is, so the communication of multiple contents can affect the strength of our ascriptions of insincerity.

2.5 That's bullshit!

In 1986, Harry Frankfurt published an influential article that identified a form of insincere discourse that didn't fit standard categories: *bullshit*. While lying requires saying what you believe to be false, bullshitting is a matter of asserting propositions whose veracity you have not even assessed. The goal of the bullshitter is often not persuasion, but some other goal, like impressing or appeasing the audience. Here's an example:

PRESIDENTIAL BULLSHIT

Pressed by journalists while leaving a government building, the president doesn't fully grasp the question she was asked. She has no idea of what the journalists are talking about, but wants to give the impression that she has everything under control. She says:

(6) The party is taking serious measures to resolve the problems you just mentioned. We have our best people working on this issue!

Here the president responds with (6) to project an image of confidence to the press. She has no idea whether the party is taking measures to resolve the problems mentioned by the journalists. As Frankfurt (1986) would put it,

[her] statement is grounded neither in a belief that it is true nor, as a lie must be, in a belief that it is not true. It is just this lack of connection to a concern with truth—this indifference to how things really are—that I regard as of the essence of bullshit. (Frankfurt 1986)

Frankfurt rightly notes that, far from being a fringe phenomenon, bullshit is widespread, notably in political speech and advertisement. The typical example of a bullshitter is the car salesman, or the politician who “never yet considered whether any proposition were true or false, but whether it were convenient for the present minute or company to affirm or deny it” (Swift 1710)³⁸.

³⁸ Sceptic philosophers, who engage in willing suspension of belief (*epoché*), famously do not have beliefs (or claim not to) on a wide variety of topics. It's unclear whether they are lying when they claim that they have no opinion about such mundane matters or whether they are bullshitting when

So far, discussion in this book has mostly assumed that either a speaker believes a proposition to be true, or they believe it to be false. Call this useful simplification the *dichotomous model* of belief, since it assumes that there are only two doxastic states: belief in the truth of a proposition, or belief in its falsity.

Bullshit throws a third state into the mix: the state of being agnostic, or *lacking* a belief, which is characteristic of this form of speech. It invites us to consider a *trichotomous model* of belief with three possible states: (i) believing p true, (ii) believing p false, and (iii) lacking a belief. These three states broadly correspond to (i) standard sincerity, (ii) standard insincerity, and (iii) bullshit.

A good theory of sincerity should presumably acknowledge that bullshitting is a form of insincerity³⁹. Intuitively, the BULLSHITTING PRESIDENT is insincere in claiming that the party is addressing the problem. However, traditional views (at least in their standard, bound formulations) deem bullshit neither sincere nor insincere. This is because bullshitters neither believe the asserted proposition to be *false* nor *true*⁴⁰.

Addressing this limitation requires integrating a trichotomous view of belief into standard accounts. This is as simple as revising the scope of the negation in any given account of insincerity – switching any requirement that “S believes *not p*” into the requirement that “S does not believe that p ”. To illustrate, here’s two “wide scope” versions of classic accounts:

STANDARD EXPRESSION (wide scope)

(inconsistently with their doctrine) they make assertions about mundane matters they are supposed to lack opinions about.

³⁹ Even broader definitions of bullshit are available, but presumably they go beyond *insincere* bullshitting, which is my focus here. Carson (2010, 62) and Stokke (2018, 140–59) argue that *evading questions* and *filibustering* with believed-true statements amounts to bullshitting. Insofar as these statements (and their implicatures) are believed to be true, they are sincere (regardless of whether they are bullshit), and so irrelevant for my purposes (incidentally, however, if this observation is correct, it undermines Stokke’s definition of insincerity, cf. Stokke 2018, 192). Cohen (2002) famously suggests that bullshit can also be a matter of content. *Unclarifiable unclarity*, or *pseudo-profound bullshit* like “hidden meaning transforms unparalleled abstract beauty” (Pennycook et al. 2015) are bullshit regardless of the speaker’s attitude towards them. I agree, but I would insist that such bullshit is insincere only when the speaker lacks a belief in its content: if they genuinely believe their own bullshit, I see no problem in saying that they are sincerely asserting pseudo-profound bullshit.

⁴⁰ By contrast, unbound views like DECEPTIVE HEARER-CENTRED classify bullshit as insincere, because of its misleadingness. As Frankfurt (1986) notes, the bullshitter “necessarily attempt[s] to deceive us about his enterprise” (i.e. whether they evaluated the veracity of their assertion, and established that it’s true).

- *A statement is insincere iff the speaker S expresses a proposition they do not believe to be true*

STANDARD HEARER-CENTRED (wide scope)

- *A statement is insincere iff it's intended by S to cause A to believe [that S believes] p, and S doesn't believe p to be true*

These wide-scope definitions of insincerity correctly classify instances of bullshit like (6) as insincere⁴¹. Since the president lacks an opinion about their statement, she doesn't believe that the content is true: both views render the verdict that (6) is insincere.

2.6 Degrees of sincerity

The opposite of truth has many shapes,
and an indefinite field
Montaigne, *Essays*, 1, IX (*On Lying*)

The trichotomous conception of doxastic states is better suited to model sincerity, especially in relation to bullshit. While this is surely an improvement over the dichotomous paradigm, it still relies on a heavily simplified picture. Belief is still modelled as a psychological state that doesn't admit of degrees. *If* you believe something, you either believe that its content is true, or that it's false – no intermediate states are allowed. Arguably, however, beliefs (and sincerity) admit various intermediate degrees; if this is right, the trichotomous model still lacks nuance.

This section will explore two ways in which sincerity comes in degrees, in ways that aren't quantifiable by trichotomous accounts. First, there are *degrees of precision*, as in statements that are believed to be *only partially true*. Second, there are *degrees of confidence*, as in statements that are *only partially believe* to be true⁴².

⁴¹ Still, hearer-centred view will struggle to accommodate the bullshit equivalent of the counterexamples discussed in §2.3 – i.e. bullshit that is not meant to persuade the audience.

⁴² The sincerity of a statement is also affected by the *degrees of strength* of the statement that conveys its content. While I have no space to discuss this here, the reader can refer to Marsili (2014; 2018a) for discussion.

2.6.1 Degrees of precision

You believe that our common friend Vladimir, who stands 190 cm high, is tall. But you have recently bought a shrinking ray. Let's remove a few centimetres from Vladimir's height, then. Now Vlad is 180 cm. Now he's 170. Now 160. As the literature on soritical progressions (e.g. Sorensen 2023) stresses *ad nauseam*, there presumably isn't a precise threshold at which our increasingly smaller friend stops being tall and starts being short (not tall). Instead, as we proceed down the line, describing Vladimir as "tall" becomes progressively less accurate. Some would say that the proposition *that Vlad is tall* becomes "less and less true".



Figure 3: A scale of shrinking Vladimirs

To explain what goes on in our mind as we observe our shrinking friend, we need to introduce some intermediate doxastic states. As we see Vladimir shrinking, we will go through a series of intermediate states not fully captured by the trichotomous and dichotomous model. After Vlad's compression begins, at some point we will opine that the proposition *that Vladimir is tall* is only "partially true", or not fully accurate, or imprecise; as the progression continues, we will eventually be inclined to think that it's false that he is tall. With some idealization, we might say that there are infinite intermediate mental states between belief and disbelief, depending on "how true" we take the proposition to be. Let's call these intermediate states *fuzzy beliefs*: they differ from unreserved full beliefs, because the observer only takes the content to be *partially accurate* (or inaccurate).

"Fuzzy logics", which allows for graded truth values, offers the resources to formalise this idea. Numerical values going from 0 (false) to 1 (true) can be assigned to propositions to measure "how true" they are. Let's then assign decreasing

numerical truth-values to the proposition that Vladimir is tall, as we evaluate it at different times: the proposition that Vladimir is tall at time 1 (call it V_1) will have a truth-value of 1, the proposition V_2 (that Vladimir is tall at time 2) a truth-value of 0.9, V_3 a truth-value of 0.8, and so forth. Intermediate degrees of truth are modelled by relying on an infinitely divisible scale of decimal numbers.

Sincerity can accordingly be conceived as a graded property, that goes from full sincerity to full insincerity, through a series of intermediate cases. Affirming that “Vladimir is tall” at different stages of the transition (at V_1 , V_2 , V_3 , etc.) is progressively less sincere. The first assertion is backed up by a belief that p is fully true (sincere statement, $V_1=1$); as we progress, this leaves ground to a belief that p is partially true (partially sincere statement, $V_3=0.8$), until we get to a belief that p is fully false (insincere statement, $V_n=0$). The example suggests that sincerity comes in degrees: some statements are more sincere than others, some more insincere. If I tell you that Vladimir is tall at V_1 , I am definitely sincere; if I tell you the same it at V_n , I’m definitely insincere. But what should be said about the cases that fall around the midpoint?

With its numerical, quantifiable account of how “true” statements are, this model offers an apparently appealing solution to this problem. There is an intermediate threshold between truth and falsity: 0.5. Statements that fall below this value are deemed *false* rather than *true*, those who fall above it are deemed *true* rather than *false*. A natural conclusion is that sincerity and insincerity, too, will follow a similar pattern⁴³. Applying this insight to the standard view, we can derive this criterion (cf. Marsili 2014, 157-8):

FUZZY

- *A speaker S is sincere iff in making their statement, S express a proposition that S believes to be more than 0.5-true*
- *S is insincere otherwise*

All idealisations, however, come at the price of some distortions. Two are worth noting here. First, as it’s currently formulated, this view posits a sudden transition from sincerity to insincerity: as soon as the 0.5 threshold is passed, a statement switches from being sincere to insincere. But as the truth value increases, we should rather expect a progressive, gradual transition from insincerity to sincerity: like “tall”, “sincerity” presumably behaves like a graded predicate with fuzzy boundaries. Ideally,

⁴³ What about statements that fall exactly in-between? I will soon argue that we should refrain from forcing intermediate cases under a specific category. According to this first pass at a definition, however, sincerity requires a truth-value *above* 0.5, meaning that a value of 0.5 yields an insincere statement.

then, an account of sincerity should acknowledge that sincerity is a graded phenomenon, which admits of degrees and lacks a sharp transition point (cf. Isenberg, 1964, p. 470; Marsili, 2014, 2018, 2022).

Second, it's psychologically unrealistic to assume that speakers invariably compute and assign a numerical truth values to propositions – especially given that there are statements that won't easily admit such numerical quantification. While partial beliefs are widespread, numerically quantifiable ones are not. To see this, consider the proposition (*p*) *that Ugo is being unfaithful to his wife*. This proposition might be “partially true” in different ways:

- (a) Ugo has let someone steal a drunken kiss from him
- (b) Ugo has a crush he will not ever act upon, but won't admit it to his wife
- (c) Ugo is flirting on dating apps, but he's not planning to meet anyone
- (d) Ugo has slept with someone else once and deeply regrets it

Under all these circumstances, *p* is at most *partially true*. But it's unclear that we can numerically quantify “how true” *p* is under each circumstance. Surely we don't reason numerically when we consider whether asserting *p* would be appropriate: there are qualitative considerations that don't easily translate into quantitative assessments. On the other hand, we seem to have opinions about how *close* to truth these statements are. Presumably, people will have intuitions about the differential closeness to truth of each circumstance. For example, one might feel that *p* is increasingly true as we move from (a) to (d) (or that some other order holds).

This suggests that sincerity has less to do with a precise numerical assessment than with some coarse-grained internal phenomenology – the speaker's perception of how true the statement is. Supposing that sincerity depends on this “perceived closeness to truth” (and building upon STANDARD EXPRESSION for simplicity), we get:

PERCEIVED ALETHIC PROXIMITY

The sincerity of a statement depends on how close to truth that statement is perceived to be by the speaker S

- *S is sincere to the extent that in making their statement, S expresses a proposition that S believes to be closer to truth than it is to falsity*
- *S is insincere to the extent that in making their statement, S expresses a proposition that S believes to be closer to falsity than to truth*

PERCEIVED ALETHIC PROXIMITY (perceived closeness to truth) is just as principled as FUZZY. Whenever the speaker numerically assesses the truth-value of the asserted proposition, these criteria output identical classifications: a statement

which is believed to have a truth value superior (or inferior) to 0.5 is also a statement that is perceived to be closer to (or further from) truth than falsity. However, PERCEIVED ALETHIC PROXIMITY (henceforth, ALETHIC PROXIMITY for short) has wider reach: it covers also circumstances where it's unrealistic to demand a numerical assessment of the proposition.

Additionally, ALETHIC PROXIMITY doesn't posit a sharp transition line from sincerity to insincerity. Instead, sincerity is presented as a matter of degree: the closer to truth a statement is, the greater its sincerity; the opposite applies to insincerity. This allows for intermediate cases: if a speaker finds it hard to establish whether the statement is closer to truth or falsity, it's equally unclear whether the statement is sincere or insincere. Rather than a sharp transition from sincerity and sincerity, we have a graded boundary and intermediate cases.

Finally, this revised criterion performs better in determining how insincere a statement is. Consider the following scenario (from Pepp, n.d.)

EXAM UNDERSTATEMENT

Two students, Jane and Sue, have just received their scores on an important exam. Jane says to Sue, "Oh Sue, I'm so upset—I only scored 55 on the exam! What did you get?" Sue looks down at her exam with a perfect score of 100 noted at the top. She knows from past experience that telling Jane the truth will only make her angry and resentful. She decides to lie about her score. Now consider two lies Sue might tell:

- (7) *I got 55*
- (8) *I got 90*

While (8) is comparatively closer to truth, it's not "partially true". On the contrary, both (7) and (8) are plainly false: they both have a truth value of 0. According to FUZZY, Sue is equally insincere in both cases. ALETHIC PROXIMITY offers a more nuanced verdict. By this criterion, both statements are insincere rather than sincere, since (from Sue's perspective) both (7) and (8) are closer to being true than they are to being false. However, (8) is correctly deemed *more insincere*, since ALETHIC PROXIMITY's general criterion states that sincerity "depends on how close to truth that statement is perceived to be by the speaker", so (7) is sincerer than (8), since Sue perceives (7) to be closer to truth than (8)⁴⁴.

⁴⁴ Pepp (n.d.) proposes a different solution. On her view, the sincerity of a statement depends both on alethic proximity and "subjective completeness" (how close a statement is to providing a complete answer to a question). This has odd consequences: uninformative tautologies (like "either

Crucially, graded insincerity isn't just a philosophical construction: it has important real-life implications. In *The Art of the Deal*, whose candid exposition of strategic dishonesty would prove remarkably prophetic, Trump writes:

The final key to the way I promote is bravado. I play to people's fantasies. People may not always think big themselves, but they can still get very excited by those who do. That's why a little hyperbole never hurts. People want to believe that something is the biggest and the greatest and the most spectacular. I call it truthful hyperbole. It's an innocent form of exaggeration—and a very effective form of promotion.

We might disagree with Trump on whether purposeful exaggeration, especially in public discourse or business negotiations, is “innocent”. But Trump's comments, if only accidentally, transpire something true – namely, that dishonest communicators like him can willingly exploit sincerity's gradedness to mitigate the social consequences of deception. If I assert propositions that are totally false, the reputational risks are high. If I make assertions that are only partially insincere, the risks are smaller: charitable audiences might think that I might just made a small mistake, or engaged in an “innocent exaggeration”. This strategy is especially effective in political discourse, where partisan affiliation might easily induce charitable interpretations⁴⁵. Studying degrees of insincerity, then, can help us better understand the grey area in which liars strive.

2.6.2 Degrees of confidence

Statements can be regarded as “truer” or “falsier”, generating intermediate doxastic states. Certainty and uncertainty, too, generate a doxastic continuum that isn't captured by the dichotomic or trichotomous model. Not all beliefs are equally strong: some we hold with more confidence (I'm certain that I'm in Madrid now), some with less (I'm somewhat confident that it's not 8 pm yet). By presenting ourselves as being more (or less) confident than what we actually are, we can misrepresent our doxastic states. Does this constitute a form of insincerity? And if so, under which circumstances?

planets have circular orbits, or they don't”, which have a minimum score for subjective completeness and maximum score for alethic proximity) turn out to be less sincere than distorting but informative approximations (like “planets describe a circular orbit around the Sun”). Intuitively, however, the first statement is fully sincere (with respect to its literal content), unlike the latter (that partially misrepresent the facts). While this view captures an interesting dimension of evaluation for discourse (e.g. the relevance of a claim, its informativeness...), it does not offer a good measure of *sincerity*.

⁴⁵ For an interesting take on the reputational dynamics of false statements directed at audiences split along partisan lines, see Saul (2024).

Like graded truth, graded confidence can be modelled by quantifying possible credal states on a numerical scale. Credences (i.e. doxastic states falling short of belief) are assigned real numbers from 0 to 1, where 0 indicates certainty in the falsity of p , 1 indicates certainty in the truth of p , and 0.5 indicates uncertainty (cases in which the subject regards p as just as likely to be true as to be false, see Figure 4).

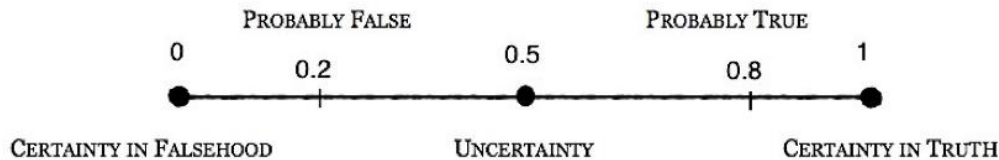


Figure 4: A visual representation of the certainty-uncertainty continuum

So, of instance, I am absolutely certain *that I am in Madrid right now* – a credence of 1. While I believe that *I will have a beer with my friends this Saturday*, I am aware that the plan might be cancelled for whatever reason – I hold a credence of 0.8. Finally, while I am convinced that (iii) *Barcelona will win the Clásico tonight*, I merely regard this as a likely possibility – a credence of 0.6.

Such graded beliefs pose challenges similar to the ones raised by graded truth-values. To illustrate, suppose that Trump, at the end of his term, exercises his bravado by claiming (9), in circumstances in which he regards this statement as *probably false*:

(9) We enacted the biggest tax cuts in American history⁴⁶

Would (9) be sincere or insincere? More generally, where should the boundary between sincere and insincere utterances be drawn when graded confidence is involved?

Once again, the temptation is to think that the threshold must fall around the midpoint, at 0.5. While this idealisation yields the right result for Trump’s (9) (classifying it as a lie), it falls into the same trap faced by FUZZY. Assuming that our credal states always admit a numerical quantification accessible to the speaker is psychologically unrealistic. When I consider whether Barcelona will win the Clásico, I might have a *feeling* of how confident I am of its victory. Unless I am in the business of gambling, however, I won’t have computed a numerical assessment of the

⁴⁶ Trump actually made this claim during his State of the Union address on January 30, 2018. Fact-checkers agree that the claim is factually false: the cuts ranked between 4th and 8th biggest in American history, depending on the calculation method.

likelihood that Barcelona will win. An account that rather relies on the phenomenology of our doxastic states (how confident *one feels*) would be preferable.

In previous work on the concept of lying (Marsili 2014; 2018a; 2023), I suggested that the sincerity of a statement depends on the speaker's perception of the likelihood of their statement. If the speaker consciously regards the asserted proposition to be more likely to be true than false, then their statement leans towards sincerity rather than insincerity; if they regard it to be more likely to be true than false, the opposite applies:

COMPARATIVE:

The sincerity of a statement depends on its perceived likelihood

- *S's statement is sincere to the extent that S takes themselves to be more confident in the truth of p than in its falsity*
- *S's statement is insincere to the extent that S takes themselves to be more confident in the falsity of p than in its truth*

Trump's assertion (9) that his tax cuts were the biggest in history is insincere by this criterion (since, *ex hypothesi*, he regards it as more likely to be false than true) even if Trump isn't fully confident in its falsity. Additionally, had Trump uttered the same statement while being *certain* of its falsity, his utterance would have been *more* insincere, since degrees of sincerity depend on the perceived likelihood of the statement. Sincere statements, too, admit similar rankings: a statement stemming from certitude will, *ceteris paribus*⁴⁷, be more sincere than the same statement stemming from lower degrees of confidence.

However, *comparative insincerity* doesn't handle the EXAM UNDERSTATEMENT example correctly. It should acknowledge that a wild departure from truth like (7) (claiming I got 55 instead of 100) is *more insincere* than a smaller departure like (8) (90 instead of 100). For COMPARATIVE, however, (7) is just as insincere as (8). Since Jane knows that she got 100, she takes herself to be maximally confident in the falsity of both propositions: both statements are considered equally (maximally) insincere.

Luckily, while COMPARATIVE cannot do what the ALETHIC PROXIMITY does (capturing closeness to truth), the opposite is arguably true: degrees of confidence

⁴⁷ This qualification is important for two reasons. First, context can affect the degree of confidence that is conveyed by a statement: depending on what's at stake (is it a matter of life and death, or a futile remark?), mutual expectations, and so forth, the very same statement may express a higher or lower degree of confidence in the same proposition (cf. Marsili 2014, 168–70). Second, there can be countervailing normative factors (like expectations of politeness) that might alter our mutual expectations of sincerity – I come back on this in §3.2.

are captured by ALETHIC PROXIMITY. If I regard a statement p to be more likely to be true than q , p is closer to truth than q from my perspective. If “closeness to truth” is interpreted in this way, ALETHIC PROXIMITY represents a better criterion, since it’s able to track two dimensions along which sincerity can be graded: *precision* (“how true” a statement is perceived to be) and *confidence* (“how likely to be true” a statement is perceived to be)⁴⁸.

2.6.3 Degrees of epistemic damage

For sympathisers of hearer-centred conceptions of sincerity, the picture sketched so far is incomplete at most: we I am yet yet to tell a story about how of confidence and perceived alethic proximity influence the hearer’s side of the picture. Krauss’s (2017) notion of *expected epistemic damage* provides an elegant remedy to this gap.

Credences can be more or less *accurate*, depending on whether they match the optimal degree of confidence that an agent should have. For example, if there’s an 80% chance that it will rain in Madrid today, a credence of 0.8 in this proposition will be maximally accurate. But testimony (what people tell us) can affect our credences: credences might become *more* or *less* accurate (or stay the same), depending on whether, by trusting what we are told, we are pushed closer or further from accuracy. Krauss calls *epistemic damage* the negative effect of testimony can have on the credences of a recipient. Epistemic damage is quantifiable: it’s the difference between how accurate the recipient’s credences are before and after trusting a statement. For example, suppose that John’s credence in the proposition that it will rain in Madrid is 0.9, but goes down to 0.3 after watching the weather report. John’s credence was only a decimal point away from maximal accuracy, but that number has now gone up to 0.5. The epistemic damage is therefore 0.4, since John’s credence has been pushed 0.4 points away from maximal accuracy.

Krauss argues that a speaker lies iff they expect to cause epistemic damage greater than 0 with respect to the proposition they asserted, conditional on the audience trusting them with respect to p . Applied to sincerity:

WORSE-OFF

The sincerity of a statement depends on its expected epistemic damage

⁴⁸ A limitation of both views is that they render unclear verdicts about cases of absolute uncertainty and bullshit. Presumably, these intermediate cases are insincere. This is an important limitation, which constitutes a trade-off for the increased nuance offered by these views (cf. §2.7.3 on trade-offs). To address it, one can opt for unbound formulations (cf. §2.4, along the lines of footnote XXX).

- *A statement with content p is sincere iff its expected epistemic damage is 0*
- *The statement is insincere otherwise*

WORSE-OFF tracks the extent to which a speaker’s deceptive intent can be graded, providing a nuanced alternative to hearer-centred views⁴⁹. However, WORSE-OFF faces substantial difficulties. First, it once again assumes that speakers compute the epistemic damage of their statements. While this is a useful idealisation, it lacks psychological plausibility – a problem shared with other views. Second, WORSE-OFF entails that a statement p cannot be insincere if it is addressed to an audience that is already maximally convinced of p . This opens the way for counterexamples, like the following (inspired by Benton 2018).

COOKIES

Kermit tells the Cookie Monster (10), while being maximally certain that (10) is false, and while being aware that the Cookie Monster is already maximally convinced that (10) is false (Elmo knows that the Cookie Monster is certain that there are some cookies left).

(10) *There are no cookies in the jar*

While Kermit is clearly insincere, he isn’t attempting to modify Elmo’s confidence in (10). The expected epistemic damage is therefore not greater than 0: WORSE-OFF classifies (10) as sincere. Generalising, counterexamples of this sort⁵⁰ suggest that there is more to insincerity than altering the accuracy of the audience’s belief. Causally contributing to preserving an inaccurate belief (as in Elmo’s example) can also constitute a form of deception (Chisholm and Feehan 1977, 144). Ideally, hearer-centred views should acknowledge that such deceptions are insincere⁵¹. Finally, WORSE-OFF remains open to many of the objections raised in §2.3 against hearer-centred accounts.

That noted, the positive observations made for hearer-centred views apply here too. *Ceteris paribus*, a statement is more insincere if both speaker-centred and hearer-centred conditions are satisfied, and the more so if expected epistemic damage is

⁴⁹ Although Krauss conceived this criterion to track degrees of confidence rather than closeness to truth, there is at least potential to extend the notion of epistemic damage to alethic proximity.

⁵⁰ See Benton (2018) and Marsili (2022) for further criticisms.

⁵¹ Standard hearer-centred views also face this objection, but can be amended to require that speaker intend to cause the hearer to believe *or continue to believe* the target proposition. WORSE-OFF, by contrast, cannot be altered in this way.

greater. Hence, like other hearer-centred views, WORSE OFF identifies an additional dimension along which insincerity can increase and decrease.

2.7 Sincerity beyond assertion

2.7.1 Beyond belief expression

Traditionally, analytic philosophers of language mostly focused on descriptive uses of speech: assertions and their truth-conditions. Other uses (such as orders, greetings, or questions) were mostly ignored. Until some philosophers reacted: Austin (1975) famously denounced this narrow focus as a “descriptive fallacy”. A new research programme, speech act theory, was born. Its main goal: to study how speech extends beyond *fact-stating* uses – how speech can perform actions (like asking, greeting, or firing someone) that go beyond describing the world.

So far, this book has knowingly committed the descriptive fallacy denounced by Austin: only the conditions under which *statements* are sincere have been discussed, limiting focus on fact-stating speech. This blind spot, however, has been strategic rather than oblivious: I am moving from simple, idealised views to progressively more elaborate ones, adding layers of complexity one at a time. This section extends the analysis to non-assertoric speech.

Philosophers often use the term “direction of fit” to distinguish two possible relations between a representation (mental or linguistic) and the corresponding state of affairs. We talk of a *fact-stating*⁵², direction of fit when the representation is meant to match the world, as it happens with beliefs and assertion. If I believe or assert “the window is open”, my belief or assertion is correct or successful if it matches this description. When instead it is the world that has to be rearranged to match a representation, we have a *telic*, or *behaviour-directing*, direction of fit. This is the case for orders and desires. If I sternly command “Open the window!”, my order (and desire) is satisfied only when facts are rearranged to match the representation (i.e. the window is open).

While the *fact-stating* vs *behaviour-directing* dichotomy captures a crucial difference, it doesn’t come close to covering the full richness of our illocutionary repertoire (the range of things we can do with language). Fact-stating uses aren’t exhausted by assertions. Some speech acts, like guesses and conjectures, present facts less forcefully than assertions; others, like solemn oaths, with even more force. Within

⁵² Here I’m using my own terminology, instead of Searle’s (1979) classic, but lexically confusing distinction between word-to-world/world-to-word directions of fit.

behaviour-directing speech, some speech acts define actions that *the speaker* has to perform (“commissives” like promises), while others acts that *the hearer* should perform (“directives” like orders and requests). Additionally, some speech acts challenge the simplistic dichotomy between *fact-stating* and *behaviour-directing*: for example, official declarations (like “I declare you husband and wife”) seem to have both directions at once.

The list could go on, but this book is about sincerity, not speech acts. For our purposes, what matters is examining how a theory of sincerity can grapple with this complexity. So, let’s see how this applies to two general conceptions of communication: the *expression view* and the *manipulation view*, and their corresponding theories of sincerity.

For the expression theorist, speech acts differ primarily because they have the function of expressing different mental states: not just *beliefs*, but also intentions, desires, regrets, and so forth. Searle (1969, 65) articulates the idea as follows:

Thus to assert, affirm, state (that p) counts as an expression of belief (that p). To request, ask, order, entreat, enjoin, pray, or command (that A be done) counts as an expression of a wish or desire (that A be done). To promise, vow, threaten or pledge (that A) counts as an expression of intention (to do A). To thank, welcome or congratulate counts as an expression of gratitude, pleasure (at H's arrival), or pleasure (at H's good fortune).

If different speech acts express different mental states, then sincerity depends not just on beliefs, but on a variety of states. Specifically, a given illocutionary act will be sincere iff the speaker has the psychological attitude expressed by that act. Conversely, it is insincere when the speaker fails to be in that mental state:

ILLOCUTIONARY EXPRESSION:

- *The performance of an illocutionary act $F(p)$ is sincere iff in uttering $F(p)$, S expresses an attitude $\Psi(p)$, and $S\Psi(p)$* ⁵³
- *The performance of an illocutionary act $F(p)$ is insincere iff in uttering $F(p)$, S expresses an attitude $\Psi(p)$, and $\sim S\Psi(p)$*

This schema yields sincerity conditions for specific speech acts, once we plug in a characterisation of a specific speech act. Paired with the orthodox assumption that assertions express belief, for example, ILLOCUTIONARY EXPRESSION yields the “wide scope” version of the standard view: an assertion is sincere when the speaker believes it, and insincere when the speaker doesn’t believe it. Plug in the standard

⁵³ $S\Psi(p)$ denotes the speaker (S) possessing the attitude Ψ (e.g. belief) with content p .

view that promises express intentions, and the result is that promises are sincere when the speaker intends to fulfil them, and insincere when they don't. And so forth for other speech acts⁵⁴.

What about hearer-centred views? For STANDARD HEARER-CENTRED, the goal of assertion is to get the audience to match the speaker's attitude: to persuade them of a content p . But STANDARD HEARER-CENTRED also admits a "higher order" variation, where the goal is to make the audience believe that the speaker has the relevant attitude (i.e. make them believe that *the speaker believes p*) as opposed to make adopt the relevant attitude (i.e. make them believe p). This higher-order variation provides is better suited to accommodate an extension non-assertoric speech. Suppose I tell you (11):

(11) I promise that I will pick up your intention at the airport

In uttering (11), my hearer-directed goal is to get you to believe that I intend to pick up your grandma at the airport: I want you to believe that I have the attitude (intention) expressed by my promise. But I don't want you to endorse the same attitude yourself (I don't want you to match my mental state, by forming an equivalent intention to pick up your grandma). The example follows the logic of higher-order HEARER-CENTRED views. Similarly for other speech acts, like thanking: if I express gratitude, it's not to get you to also feel grateful, but rather to get you to believe that I am grateful. Generalising:

ILLOCUTIONARY MANIPULATION:

- *The performance of an illocutionary act $F(p)$ is sincere iff, in uttering $F(p)$, S is attempting to induce $B(S\Psi(p))$, and $S\Psi(p)$*
- *The performance of an illocutionary act $F(p)$ is insincere iff, in uttering $F(p)$, S is attempting to induce $B(S\Psi(p))$, and $\sim S\Psi(p)$*

Like their assertoric counterparts, ILLOCUTIONARY MANIPULATION and ILLOCUTIONARY EXPRESSION often converge in their verdicts: for instance, promises will be insincere when the speaker lacks the relevant intention. However, ILLOCUTIONARY MANIPULATION requires an additional intention to get the audience to believe that the speaker is in a certain mental state. This opens the way for objections that parallel the ones reviewed in §2.3, like blatant lies.

⁵⁴ To turn ILLOCUTIONARY EXPRESSION into a view that is also sensitive to degrees of sincerity (§2.6), it can be modified so that S is insincere to the extent that S is closer to possessing the relevant attitude than they are to not possessing the attitude (conversely for sincerity). This yield verdicts comparable to *alethic proximity* and *comparative sincerity* for assertion; (cf. Marsili 2017, 148-51; 2023)

Different speech acts have different “felicity conditions”: the conditions that speakers are supposed to meet as they perform a specific speech act (Austin 1975). The views reviewed so far implicitly assume that only the violation of a subset of felicity conditions (sincerity conditions) yields insincerity. Unbound views challenge this perspective. Some philosophers (e.g. Alston 2000) argue that in performing a speech act $F(p)$ a speaker indirectly communicates that they are complying with their felicity conditions for $F(p)$. Whenever the speaker knowingly but covertly violates any felicity condition, they intentionally communicate something false – yielding insincerity on unbound views. For example, by ordering “Take a pause and drink that bottle of water!” the speaker might indirectly communicate that have the authority to issue the order, that they wish that the order to be fulfilled, that there is a bottle of water, and so forth. If any such proposition is believed to be false, the utterance can be insincere by unbound standards, even if it doesn’t violate the speech act’s felicity conditions. The net result is that any covert violation of a felicity condition count as insincerity by unbound standards.

2.7.2 Keeping your word

“Read my lips: no new taxes.”

George H.W. Bush,
(at the Republican National Convention, 8/8/1988)

During his 1988 presidential campaign, George H.W. Bush promised that “no new taxes” would be implemented. Once in office, however, he faced economic pressures, and ultimately agreed to a budget deal that included tax increases. Bush broke his promise, disappointing many Republican voters.

Promise-breaking is often considered to stand on a par with insincere assertion: we call promise-breakers liars, and both are ways to be misrepresent reality in speech. However, for every account reviewed so far (including their “illocutionary” expansions) promise-breaking is *not* a form of insincerity: insofar as the speaker *intended* to fulfil their promise at the time of making it, the promise is sincere. Sincerity is fully determined at the moment of *production* of the utterance. When one sends the message (e.g.: “no new taxes”), that message is sincere on insincere depending on one’s mental states. Whether speakers later stick to their commitments is irrelevant for traditional views.

Is orthodoxy right in maintaining that only speech production is relevant to sincerity? Or can one’s actions turn a promise uttered in good faith into an insincere promise? Arguably, “being true to one’s word” is a requirement of sincerity –

especially if we think of sincerity as a virtue, rather than a property of statements (cf. §3.2). The sincere person “talks the talks and walks the walks”, because “actions speak louder than words”.

This philosophical problem was first raised in ancient times. Assessing whether lying is always a sin, Aquinas considers the idea that “it is a lie not to fulfill what one has promised” (ST, Q110, article 3, obj. 5), only to conclude that “a man does not lie, so long as he has a mind to do what he promises, because he does not speak contrary to what he has in mind” (ST, Q110, article 3, reply to obj. 5). All the while, Aquinas acknowledges that if the promissor “does not keep his promise, he seems to act without faith in changing his mind” (ibid).

When it comes to promises that are broken for reasons that are outside the speaker’s control, it’s hard to disagree with Aquinas. Let’s go back to (11), my promise to pick up your grandma. Imagine I am intent to do it, but despite taking all precautions, my car is destroyed in an accident (for which I bear no responsibility) while I’m on my way to the airport. Surely I have been sincere, even if I broke my promise. Here the distinction between *truthfulness* and *sincerity* (§2.14) developed for assertoric discourse finds an analogue in promises. Like assertions, promises can be untruthful, or false (insofar as they fail to translate into the promised act) without being insincere.

Now consider another scenario. I promise (11), but then I change my mind. I can’t be bothered to pick up your grandma at the airport, I have better things to do – I’ll go out for a walk and an ice-cream instead. When I promised to pick her up, I genuinely intended to do it, and to stick to my promise regardless of how attractive an ice-cream-fuelled walk in the park might sound. But it turns out that I’m fickle, and I changed my mind.

Here, breaking the promise was entirely dependent on my own decisions. On top of the objective discrepancy between language and reality that is characteristic of untruthfulness, we now have the discrepancy between language and mind that is characteristic of insincerity: I expressed an intention to pick you up, but then I dropped that intention as soon as it suited me. Crucially, my change in intention constitutes a normative failure. A promissor doesn’t just communicate that they have a certain intention at the moment of promising: they guarantee that that intention is unwavering; that they will stick to it regardless of any change in their subjective preferences.

A tempting conclusion is that promise-breaking only constitutes insincerity when it’s deliberate. This “deliberate breaking” criterion distinguishes promises that are broken in good faith (as in the car crash scenario) from promises that are broken in bad faith (as in the ice-cream scenario). However, circumstances don’t always admit a straightforward under these two categories: many cases fall in-between. Sticking to a promise might impose conspicuous costs on the speaker: substantial

efforts (economic, physical), resisting social or political pressure (as for Bush's promise not to cut taxes), or genuine sacrifice (resisting to violent coercion, pain, or even death threats). Failing to resist such pressures would still constitute deliberate promise-breaking, though arguably not in bad faith. These observations only scratch the surface of the complexities we encounter in real life. It's unclear that we can draw a principled line between bad-faith cases that deserve the label of insincerity and good-faith cases that don't. Orthodoxy has the advantage of throwing such complications out of the picture: if sincere promising boils down to what the speaker intends to do at the time of the utterance, there's no need to bother with such complications.

Additionally, expressions like: "he promised sincerely, even if he ended up breaking the promise" are perfectly meaningful, even when the promise is broken deliberately. They emphasise an important distinction between honesty in speech and honesty in action. Conflating both under the same label runs the risk of impoverishing our conceptual repertoire. A preferable alternative is to use two distinct terms from what appear to be two distinct kinds of violations.

Falkenberg (1988) proposes a solution: distinguishing between insincerity and disloyalty. Disloyalty is "the intrapersonal analogue to disobedience: like commands can be disobeyed by the addressee [sic], so promises, pledges, promissory oaths or contracts can be broken by the speaker." (Falkenberg 1988: 85). Sincerity is a matter of speech production; disloyalty, a matter of speech compliance. Adopting this distinction has the virtue of preserving this book's idea that there is a core "discursive" conception of sincerity (cf. §2.1), all while acknowledging that there are other ways in which speakers can fail to communicate honestly. This richer vocabulary retains important distinctions that would be lost, were one to apply the term "insincere" to both kinds of violations.

2.8 Sincerity so far

Let's take stock. The book began by delimiting its scope: sincerity in *discourse*. This leaves aside omissions, deceptive actions, and objective falsehoods. Within the domain of communication, I identified two main families of views (speaker-centred and hearer-centred), each grounded in a different understanding of the primary function of communication (expression vs contagion/manipulation).

In response to the main difficulties of speaker-centred views (belief fragmentation and misspeaking), the *intentional expression view* offers a solution: a statement is sincere if the speaker intends to express a proposition they believe, and insincere if they intend to express a proposition they disbelieve. Hearer-centred views face more structural limitations instead, particularly in accommodating cases

where the speaker lacks any intention to persuade the audience. Nonetheless, they highlight paradigmatic cases of sincerity and insincerity, and identify one dimension along which insincerity can be incremental. The remaining sections tackled additional complications, like indirect speech, bullshit, degrees of insincerity, and non-assertoric speech. Although I've identified specific solutions for each of these challenges, I have not attempted to provide a single, unified account that handles them all at once.

There's a reason for this. Incremental adjustments are often fruitful and effective. The expression view admits an *intentional* version, which in turn admits a *wide scope* version, which in turn admits an *unbound* version: here additional refinements sum up nicely. Ideally, a good account should also accommodate non-assertoric speech (as in *illocutionary expression*), degrees of insincerity (as in *alethic proximity*), and indirect speech (an *unbound version*). While I occasionally hinted (in the footnotes) at how such incremental amendments could be achieved, I have not tried to provide an all-encompassing definition, for the result would be cumbersome and overly complex. Simplicity is also an important virtue, and whether a characterisation is satisfactory will depend in large part on what one wants to do with it: different conceptions of sincerity are appropriate for different explanatory purposes.

With this in mind, I favoured a pluralistic stance. Simple accounts, improved by small incremental amendments, can suffice in most contexts. More complex formulations can be summoned when more precision or nuance is needed. For example, when focus is on assertoric speech, operating with a definition that is sensitive to non-assertoric speech is unnecessary. But when such focus is appropriate, conceptions of sincerity apt to capture the phenomenon are available. Rather than a single correct definition, I aimed to offer a reader a diverse toolkit, which can be contextually applied, with different degrees of precision, to explore the various ways in which speech can be (or fail to be) sincere.

These different conceptions will be useful in the next section. While the first part of this book explored sincerity with a *descriptive* eye (attempting to understand what it is to be sincere and insincere), its second half covers the *normative* aspects of sincerity: what sincerity demands, and what makes it valuable.

3. The normative dimension of sincerity

I have often considered whence this custom that we so religiously observe should spring [...] that it should be the highest insult that can in words be done us to reproach us with a lie. Upon examination, I find that it is natural most to defend

the defects with which we are most tainted. It seems as if by resenting and being moved at the accusation, we in some sort acquit ourselves of the fault; though we have it in effect, we condemn it in outward appearance
Montaigne, *Essays*, 2, XVIII, (*On Calling Out Lies*)

3.1 The norm of sincerity

From a merely descriptive point of view, sincerity is just a property of statements. Understood as a norm, instead, it is a request or obligation that is placed on the act of uttering a statement. In its most general form, the norm of sincerity simply requires speakers to be sincere⁵⁵:

NORM OF SINCERITY (general)

One must: make an assertion A only if A is sincere

I've shown that sincerity can be understood in different ways. Depending on what "sincere" is taken to mean, the norm imposes different injunctions. So, for instance, by plugging the standard account of sincerity into the general formula, we obtain:

NORM OF SINCERITY (STANDARD SPEAKER-CENTRED)

One must: assert that p only if one believes that p is true

There are as many readings of the general rule as there are conceptions of sincerity. Here's a hearer-centred example:

NORM OF SINCERITY (DECEPTIVE HEARER-CENTRED)

One must: make an assertion A only if one doesn't intend A to deceive the audience

In what follows, unless otherwise specified, I will use "norm of sincerity" to refer to the general formulation, in order to keep discussion neutral and compatible with each different reading.

As for the source of sincerity's normativity, scholars disagree, along divisions that almost invariably reflect disciplinary bias. Moral philosophers tend to emphasise that sincerity is a *moral* norm, like the duty not to harm others. Linguists and

⁵⁵ The norm could also be construed as a norm demanding no insincerity (along the lines of Grice's First Maxim of Quality, cf. Grice 1989, chap. 2). Due to limitations of space, I won't discuss this variant in the book. For similar reasons, I am also limiting discussion to *assertoric* sincerity.

philosophers of language typically for granted that is a linguistic norm, and specifically a *pragmatic* norm – on the par with Gricean maxims or Austinian felicity conditions⁵⁶. Epistemologists will insist that sincerity is an *epistemic* norm; sociologists that it is a *social* norm; and so forth. Despite apparent disagreement, everyone could be right: I lean toward thinking that there's an important sense in which sincerity is a moral norm, an important sense in which it is a linguistic norm, and so forth for its epistemic or social readings⁵⁷.

Our exploration of the normative dimension of sincerity will begin by looking at the complex interactions between sincerity and other norms (§3.2). Learning how to solve normative clashes is part of what being a sincere person requires: section 4.3 provides a cursory overview of philosophical work on sincerity as a virtue. But why is sincerity deemed so valuable by philosophers? A common argument is that communication, society, and testimonial knowledge would simply not exist without a norm of sincerity (§3.4). The closing section (§3.5) will focus on its epistemic value, and on whether there are epistemic norms beyond sincerity that govern our assertoric practice.

3.2 Sincerity and other norms

Obsequium amicos, veritas odium parit
(*Adulation yields friends, sincerity enemies*)
Terenzio, *Andria* (v. 68)

La verità mi fa male, lo sai
Caterina Caselli, *La verità*

Taken in isolation, the norm of sincerity offers a recommendation that is almost comically simplistic. Being sincere is the right thing to do, and we should always be sincere. If only life was that easy! Experience teaches that reality is more complicated: our best intentions to tell the truth are often shattered by the constraints imposed on us by the circumstances of life.

Simplifying a bit, there are two main obstacles that stand in the way of our best intentions to be sincere. First, being sincere can clash with personal *interests* (what we *want* to do). In these cases, there is no question that being sincere is (*ceteris paribus*) *the right thing to do* – but in some cases, the stakes can be high, and sincerity

⁵⁶ I discussed felicity conditions in §2.7.2; Gricean Maxims are presented in the next section.

⁵⁷ Additionally, these norms could in principle come apart: for instance, moral sincerity might demand more than linguistic sincerity. I won't pursue this hypothesis here, but it's worth noting the possibility of a divergence.

truly costly. Sticking to the norm of sincerity, in this sense, is difficult because it requires sacrificing our preferences to the altar of virtue.

Second, and more interestingly, being sincere can clash with other norms (what we *ought* to do). The locus classicus is Kant's (1797) example of an axe-wielding murderer who inquires whether "our friend who is pursued by him had taken refuge in our house". Here sincerity clashes with a moral duty not to harm others – more precisely, a duty to prevent an avoidable death⁵⁸. Famously, Kant argued that our duty not to lie is perfect (i.e. exceptionless), and thus lying is never permissible – not in such extreme circumstances, and surely even less in more mundane cases, where no life is at stake.

Kant defends an "absolutist" take on the permissibility of lying. For the absolutist, no normative clash between the norm of sincerity and other norms should ever be resolved by allowing the speaker to speak insincerely: not even saving a man's life can justify lying. The absolutist view has a respectable philosophical pedigree, with sophisticated arguments put forward by influential thinkers such as Augustine, Aquinas, and Kant⁵⁹. Despite its sustained influence, however, absolutism has now gone out of fashion: most contemporary approaches allow that normative clashes can permissibly be resolved against sincerity – that is, lying is sometimes justified. Views then disagree on which considerations warrant exceptions (see Carson 2018) – a notoriously difficult task, as one moves from life-shaving scenarios to more mundane circumstances.

Clashes of sincerity with other norms happens more often than people might realise. One common clash with altruistic considerations happens when the interlocutor solicits information that would harm or hurt them. "Stick and stones will break your bones, but words will never hurt you" the refrain goes. This is certainly false, as words can be very painful – especially when they are truthful. Questions like "did you like my gift?", "do you think he/she really loves me?", or "how many days do I have left?" might elicit painful truths. Sincerity earned the nickname of "cruel virtue" (Tagliapietra 2003) precisely because being sincere, oftentimes, requires revealing to our interlocutor truths that can hurt. Navigating the narrow space between cruelty and insincerity, accordingly, is sometimes an arduous task.

⁵⁸ An analogous case (taken from Augustine) was discussed in §2.3.1: in TWO ROADS, Marcellus had to choose between an altruistic untruth (which would save his business rival) or a selfish truth (which would send him in the arms of the bandits).

⁵⁹ My succinct presentation of absolutism surely doesn't do justice to these views, which are more nuanced. For discussion, Korsgaard 1986; Sedgwick 1991; Mahon 2006; Griffiths 2004.

I mentioned how sincerity can be understood both as a moral and as a linguistic rule⁶⁰. Correspondingly, on top of clashing with moral norms, sincerity can clash with pragmatic norms. For example, Grice (1959) argues that conversations are governed by a Maxim of Quantity, which dictates:

1. Make your contribution as informative as is required (for the current purposes of the exchange).
2. Do not make your contribution more informative than is required.

Sometimes, to follow this norm (in particular its second injunction), we approximate and engage in *loose talk* at the expense of precision and sincerity. Consider the following examples (inspired by Sperber & Wilson 2002):

- (1) It's three twenty
- (2) Holland is flat
- (3) I will run to the shop before it closes
- (4) You need a Kleenex? I have one for you

It can be perfectly fine to utter (1-4) even if the speaker believes their content to be literally false. I can appropriately state (1) even if I'm aware that it's not yet 3.20 (it's 3.18), (2) even if I'm aware that Holland has some hills and depressions, (3) even if my plan is to rush to the bank without actually running, and (4) even if the tissues in my possessions are of a different brand.

According to Grice, in such cases the speaker is violating the norm of sincerity (the "Maxim of Quality", in Grice's parlance) in order to follow the Maxim of Quantity. However, unlike in genuine insincerity, where the speaker *covertly* violates the norm, here the violation is *overt* and transparent: the hearer is meant to realise that the speaker is approximating for the sake of simplicity, and that their goal is to communicate *less* than what they literally say – for example, that is *around* three twenty, or that Holland is *mostly* flat (cf. Hoek 2018)⁶¹.

It's not even clear that overt violations of sincerity like (1-4) constitute genuine violations. For speaker-centred views, the insincere speaker has to *express* the believed-false proposition, which requires representing themselves as believing its

⁶⁰ Due to limitations of space, I won't be able to discuss other common clashes, such as clashes with professional or institutional duties. Bok (1978) offers a nice discussion of these intricate issues.

⁶¹ Relevance theorists agree that this is the message to be recovered, but draw a more radical conclusion. On their view, this proves that sincerity is subordinate to another pragmatic norm – *relevance*, which demands that speakers only communicate salient content (i.e. content that is worth the audience's processing effort) (cf. Wilson 1995; Wilson and Sperber 2002).

content. But (4) presumably represents the speaker as having a tissue of some kind, rather than a Kleenex tissue specifically – to interpret the utterance otherwise would be to deliberately misunderstand it. Similarly, so long as the speaker of (4) isn't trying to convince their interlocutor that they possess a tissue of the Kleenex brand (rather than another), they won't be considered insincere by hearer-centred standards⁶².

Loose talk thus behaves like metaphors and irony: by openly violating the norm of sincerity (*flouting* it, in Grice's parlance) the speaker manifestly intends to convey a true proposition. This is exemplified by (5) and (6), which respectively communicate that the math problem is easy (rather than a pastry product) and that smoking is unhealthy⁶³:

- (5) This math problem is a piece of cake!
- (6) Smoking is *great* for your health... *of course it is*.

A similar phenomenon occurs when sincerity clashes with norms of politeness. When there's mutual awareness that politeness clashes with sincerity, this can affect what we ultimately take the speaker to communicate. Consider the following:

- (7) Beppe: How are you?
- (8) Lewis: I'm fine
- (9) Boris: Congratulations!

If Beppe and Lewis are barely acquainted business competitors, there is no expectations that Lewis respond to (7) with anything else but (8): bringing up their personal life's difficulties would be socially inappropriate in the context. Similarly, an expression like (9) is often expected regardless of whether the speaker genuinely feels happy about the achievement of the recipient. Given that politeness demands the use of these expressions even in the absence of the mental state that they supposedly express⁶⁴, we do not always take them to communicate these states. By uttering them, the speaker sometimes conforms to a social ritual (like saying "hello!" or "goodbye!") without communicating much, beyond their willingness to conform to such ritual. If this is right, *at least some* polite falsities aren't truly insincere, because

⁶² For ALETHIC PROXIMITY, (1-4) are even more straightforwardly sincere, since they don't stray too far from what the speaker considers true.

⁶³ If the speaker does not believe that the problem is easy and that smoking is unhealthy, they can still be insincere. While bound definitions will struggle to acknowledge this, unbound definitions naturally yield this verdict. This phenomenon is amply discussed in contemporary literature under the label of "non-literal lying" (Viebahn 2021; Marsili and Löhr 2022; Güngör 2024).

⁶⁴ Or that they supposedly aim to induce in the audience, if you're a manipulation-theorist.

no believed-false content is communicated. Which is not to say, of course, that all polite falsities are sincere. Consider the classic:

(10) You look *soooooo* good with your new haircut!

Presumably, polite falsities like (10) rather lean toward the insincerity end of the scale, especially if volunteered by the speaker, rather than elicited by a question. In other cases, weighing politeness against sincerity will be more difficult, yielding mixed verdicts. The main takeaway is that *some, but not all*, polite falsities can be excused because our social conventions are such that the speaker ends up communicating less than what they literally say.

3.3 Sincerity as a virtue, insincerity as a vice

I mentioned (§2.1.1) that sincerity can be understood both as a *property of utterances* (“discursive sincerity”, our focus so far) or as a *property of people* (“dispositional sincerity”). The two notions, of course, are related. A sincere person is a person who has a stable disposition to make sincere *statements*. If the virtue of sincerity amounts to having this disposition, this opens the way for different understandings of it, depending on what “sincerity” is taken to mean. For example, on a standard speaker-centred view, the virtue of sincerity is:

SPEAKER-CENTRED VIRTUE OF SINCERITY

A sincere person is one with a stable disposition to only assert what they believe to be true

What about the vice of insincerity? Despite growing interest in vice epistemology, I am not aware of any systematic discussion of this vice: insincerity is often listed next to other vices, but not discussed. Let’s try to fill this gap. It’s tempting to define the vice as the mirror image of the virtue, simply switching “true” for “false” in the schema. However, only chronic liars have a stable disposition to only assert disbelieved propositions. A more promising alternative is to understand the vice of insincerity as the absence of the virtue of sincerity. The resulting conception rings more plausible:

SPEAKER-CENTRED VICE OF INSINCERITY

An insincere person is one without a stable disposition to only assert what they believe to be true

Virtue theorists (in ethics and epistemology) sometimes emphasise that being virtuous requires valuing certain conducts (like sincerity) as intrinsically desirable. Norms like the norm of sincerity, by contrast, are often taken to be sustained by an infrastructure of social policing: speakers are motivated to be sincere because infractions are socially sanctioned. To be motivated by such instrumental considerations, some argue, falls short of possessing the virtue of sincerity. For Williams (2002, 59), the truly sincere person is disposed to be sincere not because sincerity has instrumental value (e.g. preserving one's reputation), but because they recognize sincerity as intrinsically right and desirable. The virtue of sincerity, on this view, requires being also moved by a genuine appreciation of the worth of being sincere.

Historically, Aristotle's *Nicomachean Ethics* set the ground for philosophical discussion of the virtue of sincerity. For Aristotle, the sincere or truthful person (*aletheutikos*) embodies a mean (*mesotes*) between two vicious extremes: the boastful exaggeration of the braggart and the sly dissimulation of the self-deprecator. Aristoteles's discussion focuses on discourse *about oneself*. The sincere person tells the truth about oneself, no more and no less than what one deserves. (NE, IV, 7, 1127a 17-28).

Discussing the virtue of sincerity, Aquinas (ST, II-II, 109, a1) aligns with Aristotle's position. Aquinas sees sincerity as an exclusively moral virtue, arguing that it is "neither a theological, nor an intellectual, but a moral virtue". This Aristotelian take is in explicit contrast with contemporary philosophy – in particular virtue-epistemology, where sincerity is regarded as a paradigmatic epistemic virtue.

Like Aristotle, Aquinas acknowledges that understatement is a lesser threat to sincerity than overstatement. After all, the booster generally gains undeserved credit, whereas the modest speaker avoids ostentation, as Socrates himself did (NE IV, 7, 1127b 25-26). Additionally, understating can be "done without prejudice to truth, since the lesser is contained in the greater" (ST II-II, 109, a4). Contemporary work agrees on this: as assertoric force increases (as in stronger statements or boosting), the same communicated content is *ceteris paribus* more insincere (Marsili 2014; 2018a). Excessive modesty (understatement), however, is seen as a vice by some philosophers, and receives the highest condemnation from Montaigne:

As to this new virtue of feigning and dissimulation, which is now in so great credit, I mortally hate it; and of all vices find none that evidences so much baseness and meanness of spirit. 'Tis a

cowardly and servile humour to hide and disguise a man's self under a visor, and not to dare to show himself what he is

A separate issue is whether sincerity, as a virtue, requires a disposition to *disclose* one's beliefs. Some philosophers hold the strong view that the truly virtuous speaker has to be *supersincere* (§2.1.2). Rousseau's (RSW) radical take, for example, is that sincerity requires full transparency, including volunteering information that jeopardises one's reputation and social standing⁶⁵. A less radical view is that speakers should not deliberately withhold information, and should reveal their beliefs when this is socially expected and helpful (Williams 2002, chap. 5; Queloz 2021, 165–66).

I have already stated my (mainly terminological) case against this view. If our concern is *sincerity in discourse*, sincerity cannot impose demands on what one has not yet said. Understood as a virtue exercised in discourse, then, sincerity merely requires that *when* something is communicated, it matches one's belief – as such, it doesn't forbid omission or dissimulation. To put it in Montaigne's words: "A man must not always tell all, for that were folly: but what a man says should be what he thinks, otherwise 'tis knavery".

This is not to deny that supersincerity might be demanded by other virtues. *Spontaneity* (though not necessarily a virtue) is a disposition to telling immediately and directly what you believe – a disposition to share's one thoughts in an unmediated way. *Frankness* is often understood as a disposition to disclose hard truths – even those that can hurt our interlocutor, or that might clash with social expectations. The courageous revelation of hard truths is also required by *parrhesia*, a disposition (famously discussed by Foucault 2011) to speak the truth, especially against repression and silencing by oppressive authorities. Spontaneity, frankness and parrhesia all demand some degree of supersincerity, since require voluntary disclosure of information – often, information that the speaker would profit from keeping to themselves.

Is sincerity the manifestation of some other more general virtue? For Miller (2021, 22–23) sincerity should be understood as a form of honesty⁶⁶. Honesty, in his view, requires a general disposition to act in a way that doesn't intentionally distort the facts as the agent sees them. For instance, a cheater typically misrepresents the facts about whether they are following the relevant rules. They are dishonest, but not insincere (insofar as their cheating doesn't involve communicating

⁶⁵ For a critical discussion of Rousseau's take on sincerity, see Williams (2002, ch.8).

⁶⁶ Miller uses a slightly different vocabulary: he uses "truthfulness" for the speaker-centred virtue of sincerity, "forthrightness" for its non-literal counterpart, and "veracity" for both (minor details aside). Since "veracity" and "truthfulness" are often understood to refer to the communication of *objective* truth, I have here adopted a different terminology.

insincerely). If this is right, misrepresentation of facts, which some might intuitively regard as the hallmark of insincerity, is central to all dishonest behaviour. Sincerity, in turn, might be regarded as the manifestation of honesty in speech – a disposition *to communicate* in a way that doesn't intentionally distort the facts as the agent sees them.

3.4 Sincerity: language, knowledge and society

Philosophers hold sincerity in high regard. They often depict it as a fundamental virtue, or an exceptionless norm, laying at the foundations of human morality. Why is sincerity considered so important and valuable? What purpose does it serve in our society? Let's review some answers, focusing on sincerity's role in sustaining communication (4.3.1-2), meaningful social relationships (4.3.3), and testimony (4.3.3).

3.4.1 Sincerity as the foundation of communication

Could there be a linguistic community where people are insincere virtually all the time, but still manage to communicate? Coady (1992) imagines alien world of this sort, where individuals never speak the truth – let's call this planet Lars, (like Mars, but with an L for lying) :

LARS

Let us suppose for the moment that [the Lartians] have a language which we can translate (there are difficulties in this supposition as we shall see shortly) with names for distinguishable things in their environment and suitable predicative equipment. We find however, to our astonishment, that whenever they construct sentences addressed to each other in the absence (from their vicinity) of the things designated by the names [...], then they seem to say what we (more synoptically placed) can observe to be false. (Coady 1992, 85)

Can we really say that the Lartians (the inhabitants of Lars) have a practice of *asserting*, or of *reporting* what happens in the external world? Let's assume that expressing a belief is essential to communication and assertion, in the spirit of "expression theories". Since no correlation between the speaker's assertions and their beliefs is ever observed on Lars, it's hard to see how assertions could be regarded as devices for expressing beliefs. By expressivist standards, then, no genuine assertions can be made on Lars. A similar conclusion follows from manipulation views. The function of assertions cannot be to persuade audiences on Lars: given that their statements are never observed to match reality, Lartians

speakers cannot be intending to persuade even the most gullible of audiences. Whatever Lartians are doing with declarative sentences, it's something different from our practice of asserting⁶⁷.

Widespread insincerity presumably breaks down also the smaller components of language: the meaning of nouns and predicates. If a world like “walrus” refers to a certain category of mammals, it's not because this arbitrary string of letters has some intrinsic relation with the marine species. Rather, it's because our linguistic community established a convention, and English speakers coordinate to use the term in this way. It's hard to imagine how such coordination could be achieved on LARS. In situations where a walrus is present, Lartians will indicate and proffer all sorts of expressions that don't involve the string of words “walrus”: things like “there's a boat there”, or “there's Henri Kissinger there”. *In virtue of what* would words have the meaning they have in LARS? Coady (1992, 83–89) is highly sceptical that we can come up with a plausible theory of meaning that addresses the question.

For similar reasons, language acquisition would be practically impossible on LARS. Lartian kids will hear all sort of things when they're around walruses: adults will indicate the animal and call it a boat, Henry Kissinger, and all sorts of other things. If language acquisition requires repeated exposure to regularities, and a language without sincerity breaks down those regularities, then language acquisition isn't possible – or at least highly difficult – on LARS.

Without a modicum of sincerity, no meaningful linguistic communication can occur – this, at least, is what many philosophers conclude from thought experiments like LARS⁶⁸. Some take this idea further, and argue that sincerity is the foundation of language itself. In *Languages and Language* (1975), Lewis (building on Stenius 1967) argues that truthfulness and trust are essential to linguistic communication. For Lewis, to have a language at all just is to have such a convention⁶⁹:

My proposal is that the convention whereby a population P uses a language \mathcal{L} is a convention of truthfulness and trust in \mathcal{L} . To be truthful in \mathcal{L} is to act in a certain way: to try never to utter any sentences of \mathcal{L} that are not true in \mathcal{L} . Thus *it is to avoid uttering any sentence of \mathcal{L} , unless one believes it to be true in \mathcal{L}* . To be trusting in \mathcal{L} is to form beliefs in a certain way: to impute truthfulness in \mathcal{L} to

⁶⁷ The same conclusion is inescapable even for the orthodox view that asserting requires *presenting a proposition as true* (a requirement that is rarely questioned in the literature – cf. Marsili and Green 2021, 23–26). Since no correlation can ever be observed between statements and reality on Lars, Lartians could hardly see themselves as being in the business of presenting the content of their utterances as descriptions of the world.

⁶⁸ For a different take, see e.g. Skyrms (2010, 73–82).

⁶⁹ Conventions are understood by Lewis as arbitrary, self-sustaining solutions to coordination problems (for a full definition, see Lewis 1975, 4-7).

others, and thus to tend to respond to another's utterance of any sentence of \mathcal{L} by coming to believe that the uttered sentence is true in \mathcal{L} . (*Lewis, 1975, 7, my emphasis*)

On this view, linguistic communication is made possible by a convention that speakers try to make assertions that are true (truthfulness), so that hearers generally interpret assertions as true (trust). This enables coordination between communicators. Unlike Coady's Lartians, the Lewisian "population P" experiences a regular (albeit not perfect) correspondence between what people say and states of the world, so that none of the Lartian difficulties (in language acquisition, interpretation, etc.) arises for them.

Grice (1989, 27) also regards sincerity⁷⁰ as a prerequisite for meaningful communication. In his "Retrospective Epilogue" he emphasises that unless expectations of sincerity (in his terminology, "the maxim of Quality") are in place, utterances fail to be genuinely communicative (in his terminology, "contributions"):

The maxim of Quality, enjoining the provision of contributions which are genuine rather than spurious (truthful rather than mendacious), does not seem to be just one among a number of recipes for producing contributions; it seems rather to spell out the difference between something's being and (strictly speaking) failing to be, any kind of contribution at all. (Grice, 1989, 371)

3.4.2 When insincerity spreads

Some philosophers took this line of reasoning even further. For Immanuel Kant, truthfulness is an unconditional duty that is essential not only for meaningful language use, but for the very possibility of society. Kant goes as far as suggesting that every single lie threatens the collapse of language and of institutions (Kant 1797):

[When I lie,] I cause that declarations should in general find no credence, and hence that all rights based on contracts should be void and lose their force, and this is a wrong done to mankind generally.

⁷⁰ To be precise, both Lewis's and Grice's discussion oscillates between two notions: the speaker-centred norm of sincerity, and the requirement that the speaker should *aim (or try) to tell the truth*. Lewis's definition of truthfulness appeals to the latter, but immediately specifies that it entails the former (see my emphasis in the quoted text). Grice has different labels for the former notion (the "first Maxim of Quality") and the latter (the "Supermaxim of Quality"); though ambiguous, the terminology used in the passage suggests that its focus is on the former.

Taken literally, Kant’s position strikes as bizarrely radical: how can a single lie have such profound consequences on the fabric of society? But his remarks point in the direction of two plausible claims. The first is that lies pose an incremental threat to the credibility of assertions: the more lying spreads in a population, the less credible each assertion made by an individual that belongs to that population⁷¹. The second is that sincerity matters not only for communication, but for society at large. The point was already raised in Aquinas’s *Summa Theologiae*, where the Doctor Angelicus writes that:

Since man is a social animal, one man naturally owes another whatever is necessary for the preservation of human society. Now it would be impossible for men to live together, unless they believed one another, as declaring the truth one to another. (Summa Theologiae, 1, 109, a. 3)

If Aquinas is right, a functioning society requires modicum of interpersonal trust, and that trust is difficult to achieve if one cannot expect other people to be sincere. Sincerity is a *conditio sine qua non* for human civilization. This is in line with Kant’s point that without sincerity we could not have contracts and institutions more generally. Modern societies rely on a communication and mutual trust to create complex networks of obligations. Institutions and social objects (nations, money, corporations, contracts) all rely on the presumption that people, by at large, do what they say. Without sincerity and loyalty (as defined in §2.7.2), there can be no functioning society.

3.4.3 Sincerity and knowledge

There’s another human faculty that presumably depends on sincerity: our ability to learn from one another. There is no doubt that much of human knowledge has been passed on through communication, rather than by direct acquaintance with the known facts (e.g. Coady 1992; Alfano and Levy 2020). In the late 80s, philosophers grew increasingly more interested in the importance of *testimony*. “Testimony” here designates any speech act that presents a proposition as true⁷², typically to persuade

⁷¹ This consideration has been well-explored in the game-theoretic literature, where mathematical models are adopted to study how the spread of deceptive signallers threatens to undermine a population’s ability to communicate meaningfully (e.g. Skyrms 2010). These conclusions of these studies, however, are less dire than Kant’s: a communicative system can retain stability (in game-theoretic sense) despite widespread deception, and senders can successfully convey information even when deception is universal (Skyrms 2010, 73–82).

⁷² This characterisation will inevitably find detractors, but it is sufficiently general for current purposes. On defining of testimony, see Lackey (2006) and Cullison (2010).

an interlocutor. Testimonial beliefs, in turns, are beliefs acquired by accepting testimony.

Historically, philosophers have been disagreeing about the conditions under which testimonial beliefs can be said to be *warranted* or to amount to knowledge⁷³. Few, however, would disagree that sincerity is an important prerequisite for our ability to acquire testimonial knowledge. In Coady's Lartian society, where expectations of sincerity are absent, testimonial knowledge is surely beyond reach: if Lartians were able to communicate, they certainly could not trust the reports of their fellows.

Let's grant that we can learn through testimony, and that this is at least in part because that this process of information acquisition, in combination with other cognitive faculties⁷⁴, is a sufficiently reliable source. A question remains open: which role does sincerity play in our ability to acquire testimonial knowledge? And is sincerity all we demand from testimony?

3.5 Testimony and the norm of sincerity

3.5.1 Testimony and epistemic norms

Philosophers often emphasise the importance of *epistemic norms* (and obligations) in sustaining the reliability of testimony (Williams 2002; Goldberg 2011; 2015). In a nutshell, if a good proportion of testimony is veridical, it's because speakers keep each other in check, ensuring that veridical testimony gets rewarded, and unreliable testimony punished (Marsili 2025). This (extremely simplified) story leaves open an important question. Which epistemic norm governs this process of information sharing through assertions?

Since Williamson (1996), philosophical discussion of this question has taken a simplified form. Scholars often assume that assertion is governed by a single norm of the form "assert that *p* only if *p* has C", where C is an epistemic property like *belief*, *justification*, *truth* or *knowledge*. Unlike norms of politeness, or other conversational norms (such as the Gricean maxims of Quantity discussed in §3.2),

⁷³ *Non-reductionists*, who hold that we have *pro tanto* epistemic entitlement to trust (a piece of) testimony, oppose *reductionists*, who deny that we have such epistemic entitlement (for an introduction, Leonard 2023). Crucially, both agree that trusting testimony *under the right conditions* yields knowledge.

⁷⁴ This qualification includes our ability to engage in *epistemic vigilance* (Sperber et al. 2010) – to assess sources for their reliability, and content for its plausibility. The *process* of acquiring belief through testimony is therefore understood to include epistemic vigilance, making the claim compatible with both reductionism and non-reductionism.

the norm of assertion regulates *only the speech act of assertion* – it is distinctive of it⁷⁵. If this simplifying assumption is granted⁷⁶, identifying the norm of assertion is a matter of determining which single epistemic property makes a proposition assertable. The following proposals dominate the academic debate:

(KR) KNOWLEDGE-RULE: “Assert p only if you know that p”

(TR) TRUTH-RULE: “Assert p only if p is true”

(JR) JUSTIFICATION-RULE: “Assert p only if you rationally believe that p”

(BR) BELIEF-RULE: “Assert p only if you believe that p”⁷⁷

Typically, defenders of more demanding norms (like KR, JR, or TR) acknowledge that speakers are expected to conform to a sincerity-rule like BR⁷⁸. However, there are exceptions: it has been argued that BR *doesn't* regulate assertion, and that we could count on testimony to be reliable even in its absence. The next section presents and tackles these objections. Established that BR is *necessary* for permissible assertion, the rest of this section deals with its *sufficiency*: whether the epistemic norm in force in our linguistic community is, in fact, more stringent.

⁷⁵ And therefore of testimony, if we accept the additional assumption that assertions present their content as true (cf. footnote XXX).

⁷⁶ For a review of reasons not to accept it, Marsili & Pagin (2021, Supplementary Document: Which Kind of Norm?).

⁷⁷ This list is inevitably not exhaustive, and contains important approximations. It groups together views that differ in important ways (such as *justification rules*, which differ radically from one another). It also leaves out many prominent alternative views, like the context-sensitive proposal defended by Goldberg (2015). For a more detailed review, see Pagin and Marsili 2021.

⁷⁸ Since you can only know what you believe, KR requires conformity to BR. JR is often understood by their proponents to require a reasonable or justified belief. The story is more complicated for TR, which requires belief via “secondary norms”. Secondary conformity to a norm allegedly requires that you reasonably believe that you are following the norm. TR thus “secondarily” requires the speaker reasonably believe that the proposition is true (Weiner 2005; Whiting 2012). For criticism of this move, see (Douven 2006, 478–80; Lackey 2007; Gerken 2011; Schechter 2017) and footnote XXX.

3.5.2 Against sincerity: selfless assertions

According to BR, the *belief-rule*, an assertion is epistemically permissible only if it's believed to be true. This belief-rule is effectively a *sincerity-rule*, since it says that an assertion is appropriate only if it is sincere in the speaker-centred sense⁷⁹.

Philosophers rarely question the idea that BR governs assertoric speech. After all, “don't lie” is one of the first pragmatic rules that we learn as we begin to speak. Expectations of sincerity are so deeply engrained that we tend to take them for granted, both in life and in philosophical theorising. But not all philosophers agree that they are universal.

The most prominent critic of BR is Jennifer Lackey, who argues that justification, but not belief, is required for appropriate assertion. She defends the following variant of JR:

(JR) EXTERNAL-JUSTIFICATION-RULE: “Assert p only if it is rational for you to believe that p ”*

According to JR*, an assertion is epistemically permissible only if it is rational for the speaker to believe that p , regardless of whether they actually believe that p . The underlying idea is that speakers should strive to assert not what they personally feel is true, but rather the propositions that are best supported by the evidence available to them.

Of course, JR and BR often converge in their predictions: if it is rational for me to believe that p , I will typically believe that p . But what about those cases when justification and belief come apart, and a speaker fails to believe what they are justified to believe (or vice versa)? Lackey (2007) claims that it is precisely in these cases that we can notice the superiority of JR*. Only this rule can accommodate the intuitive appropriateness of *selfless assertions*, which she defines as follows:

⁷⁹ Manipulation views, too, can be linked to BR. If assertors invite their audience to believe what they assert, then cooperation demands that you refrain from asserting unless you believe the proposition yourself. This is the familiar Gricean picture, and Bach (2008) defends BR precisely by appealing to this principle.

An assertion that p is selfless if and only if:

- 1. a subject, for purely non-epistemic reasons, doesn't believe that p;*
- 2. despite this lack of belief, the subject is aware that p is very well supported by all of the available evidence; and*
- 3. because of this, the subject asserts that p without believing that p*
(Lackey 2007, 599, substantially edited)

The most discussed example of selfless assertion involves a doxastic conflict between religious faith and scientific evidence (Lackey 2007).

CREATIONIST TEACHER

Stella, a creationist teacher, is aware that evolutionary theory is supported by the best available evidence, but firmly accept creationism (which contradicts evolutionary theory) on the basis of her religious faith. Suppose that Stella tells her students (11), which she believes to be false:

(11) *Homo Sapiens evolved from Homo Erectus*

Intuitively, it would be appropriate for Stella, as a teacher, to assert (11): after all, she is presenting what she takes to be the scientific consensus on the subject. In a classroom context, this seems appropriate. But BR incorrectly predicts that (11), a believed-false statement, should strike as inappropriate and criticisable. Generalising, the intuitive appropriateness of selfless assertion challenges the universality of expectations of sincerity, and appears to corroborate Lackey's view that speakers are only expected to make justified assertions, regardless of whether they believe them.

However, there are several considerations that should give us pause. First, it's unclear that selfless assertions are genuine assertions. Teachers are under a professional obligation to teach students the scientific consensus, not their personal opinions. Utterances like (11) could then be understood as a speech act that fall short of assertion (e.g. an act of *reporting* the theories that are presented in the textbook) or as an assertion with an implicitly prefaced content (e.g. "*According to the scientific consensus*, Homo Sapiens evolved from Homo Erectus). Under both interpretations, Stella's utterance would not be insincere, nor would it prove that there are permissible violations of BR and KR. It has been argued (Milić 2017) that this treatment extends to other examples discussed in the literature, too, since they involve statements made under professional obligation (e.g. a vaccine-sceptic *doctor*, a racist *juror*, etc.).

In response, for each scenario, one might imagine the speaker repeating the same statement in a different context – talking to a friend at the pub, for instance⁸⁰. However, once we move to non-professional contexts, intuitions about the appropriateness of the statement change. Imagine that Stella asserts (11), without any qualification, to a childhood friend in a pub. Here there’s a stronger case to be made that Stella is lying, or at least that her assertion is defective in important ways – she is, after all, falsely presenting the scientific consensus as her own opinion, without any qualification (cf. Milić 2017, 2291-93).

Relatedly, if one accepts that selfless assertions are genuine assertions, it’s unclear that they are appropriate. Selfless assertors always have the option to use fuller, more informative expressions, like:

(12) According to the available evidence, (11). But if you are asking me what I think, not (11)

By choosing to make a plain statement (*p*) instead of using more informative expressions like (12), selfless assertors misrepresent themselves as believing what they do not actually believe. They imply, or otherwise convey, that they believe what they are saying, thereby inviting their audience to form a false belief, when this is perfectly avoidable – given the availability of (12). By most unbound standards, (12) is therefore insincere. Accordingly, there is an important sense in which the assertion is defective, criticisable, or even epistemically impermissible (given the available alternatives)⁸¹.

Suppose you have the opposite intuition. Stella is aware that (11) has good epistemic credentials, and that is therefore likely to be true: intuitively, then, her assertion is permissible. Empirical evidence (Turri 2014) shows that many non-philosophers share this opinion: they consider Stella’s assertion permissible. However, the same study also found that when participants judge a selfless assertion to be permissible, they almost invariably interpret the scenario as one in which the protagonist believes that what they are saying is true. In other words, when people feel the pull of the intuition that (11) is appropriate, it’s typically because they regard Stella’s awareness of the evidence in support of (11) as an indicator that there’s an important sense in which she doesn’t believe (11) to be false. This, and not the fact

⁸⁰ This is how Lackey (2007) construes the *racist juror* scenario, anticipating this worry.

⁸¹ Ironically, the verdict that Stella’s assertion should feel defective is inescapable given Lackey’s own framework, since she argues that communication is governed by the NMNA (“S should assert that *p* in context *C* only if it is not reasonable for S to believe that the assertion that *p* will be misleading in *C*”) in addition to JR*, and violations of NMNA should be intuitively inappropriate in her view. So, Lackey either has to drop the claim that selfless assertions are intuitively defective, or the claim that NMNA governs speech in general.

that no sincerity-rule is violated, is what drives laypeople's intuitions that selfless assertions are permissible. This should give us pause before we draw strong theoretical conclusion from intuitions about these cases.

In conclusion, selfless assertions are too controversial to serve as decisive counterexamples. They're complex examples, whose interpretation is debatable. It's not clear that they are genuine assertions. It's not clear that they are intuitively appropriate. And it's not clear that our intuitions about them aren't polluted by some natural propensity to ascribe beliefs to the speaker as if they were thinking rationally. Before abandoning the strongly intuitive thesis that assertions are governed by a sincerity-rule, convincing evidence is needed. While selfless assertions complicate the picture, they do not provide the knockout argument that would be needed to abandon the consensus view that there is a *pro tanto* expectation that people believe what they assert⁸².

3.5.3 Unjustified beliefs and assertions

To accept that assertions are governed by BR is not yet to agree that BR is the *only* epistemic requirement for proper assertion. In fact, most philosophers reject this thesis, since assertions can be epistemically faulty even when they are sincere. Sometimes people form beliefs on very poor grounds – a hunch, wishful thinking, fallacious inferences, etc. Let's call such poorly formed beliefs "unjustified beliefs", to emphasize that they are not supported by appropriate evidence or reasons. When speakers assert such beliefs without qualification, they make *unjustified assertions*. In most contexts, such assertions are epistemically faulty. Consider an example:

GOOD COOK

Bob and Rachel are having a conversation about their common friend Jacques. Rachel asks whether Jacques is a good cook. Bob does not know. However, Bob knows that Jacques is French, and he is under the impression that French people are often excellent cooks. On this basis, he replies:

(13) *Sure, Jacques is an excellent cook*

In asserting (13), Bob makes an *unjustified assertion*. Intuitively, (13) is epistemically inappropriate⁸³ – presumably, because Bob lacks sufficient ground to make his claim.

⁸² For some additional considerations against BR, see Wilson 1995; Wilson and Sperber 2002; Mandelkern and Dorst 2022.

⁸³ The intuition that UNJUSTIFIED ASSERTIONS are inappropriate is strong, cross-cultural, and widespread. For a review, see Graham and Pedersen (2024) and Kneer and Marsili (2025).

Imagine that you later realise that Bob had no evidence for (13) beyond Jacques' nationality. It would be natural (and appropriate) for you to complain, or to criticise Bob for his unfounded claim. More generally: unjustified assertions are inappropriate *qua assertions*, and BR alone cannot accommodate this fact. This fact is equally taken for granted in post-Gricean pragmatics, since Grice's (1989) Second Maxim of Quality dictates that cooperative communication requires appropriate evidence in addition to sincerity (cf. Searle 1969, 66). Philosophers who argue that assertion is governed by a stricter rule (like JR or KR), too, underscore the impermissibility of unjustified assertions.

3.5.4 The compound view

The shortcomings of BR can be addressed by pairing BR with a norm that forbids the acceptance of unjustified beliefs – this, at least, is the solution propounded by the “compound view”. This view comes in two flavours, depending on whether the compounding requirement is understood as a *norm of belief* or as an *epistemic virtue*. Let's start with the first.

Bach (2008) and Hindriks (2007) suggest that while assertion is only subject to BR, beliefs are themselves subject to norms – specifically, a norm that demands that one should believe only what one knows⁸⁴. Once this requirement is paired with BR, we get the derivative requirement that one should only assert what one knows. Hindriks (2007, 23) resumes the derivation as follows⁸⁵:

[BR] One must: assert that *p* only if one believes that *p*.
 [KR-B] One must: believe that *p* only if one knows that *p*.
 ∴[KR] One must: assert that *p* only if one knows that *p*.

If KR can be derived from BR, unjustified assertions are forbidden, reconciling BR with our intuitions about which assertions are permissible. Before Bach and Hindriks, Williams (2002) has defended another version of the “compound view”, which compounds BR with two “virtues of truth” rather than norms. These virtues are *Sincerity* and *Accuracy*. I already discussed the former (§3.3). Accuracy, instead, involves (broadly) a disposition to form beliefs in a careful way. A careful believer displays a sensitivity to evidence, a concern for truth, and a disposition to “get things right”. Like Bach's or Hindrik's knowledge norm of belief, it ensures that people don't form unjustified beliefs. But there are key differences.

⁸⁴ Whether this claim is plausible is itself matter of controversy (McGlynn 2014; Hughes 2017).
⁸⁵ I altered the acronyms to match the ones I adopt. Additionally, Hindriks qualifies BR and KR so that they hold only “in situations of normal trust”.

First, Williams’s emphasis is on accuracy (and therefore a stable correlation with truth), rather than knowledge. Second, Williams believes that only a virtue that is valued in its own right (rather than a norm that is followed for instrumental reasons) can sustain testimonial knowledge. Despite these differences, the virtue-theoretic compound view yields a similar picture. The virtue of sincerity ensures that what people assert, by at large, tracks what they believe. In turn, accuracy ensures that what people believe, by at large, tracks the truth. The result is a community in which assertions reliably convey beliefs that are true, and therefore in which we can trust that testimony is, by at large, true.

Compound views postulate a normative “division of labour” between doxastic and assertoric normativity. An *assertoric* norm (or virtue) of sincerity ensures that our claims reliably express beliefs; a *doxastic* norm (or virtue) of accuracy for belief ensures that our beliefs reliably track the truth. On this conception, sincerity is all that assertion demands: BR exhausts assertoric normativity.

However, compound view faces compelling objections. First, there is the *derivation objection*. Both norms of assertions and norms of beliefs are often described as “epistemic norms”. Their function (e.g. maximizing knowledge or truth) and content (belief is a condition for knowledge) certainly warrants such terminology. Nonetheless, there are important senses in which these norms are structurally different. For example, it could be argued that their normative force is *grounded* in different facts. Many philosophers recognise that norms of assertions are grounded in contingent social facts: whether a certain norm governs assertion depends on whether it is *in force* in a certain community (García-Carpintero 2022) that polices its infractions (Alston 2000, 251–62). The same is certainly not true of norms of belief.

Given that, despite the common label of “epistemic”, there are asymmetries between the norms, it’s unclear that we are entitled to derive KR from KR-B and BR. If assertion requires belief (because of our socio-linguistic practices), it doesn’t follow that it also requires an *epistemically responsible* belief (unless epistemically responsible belief is what our socio-linguistic practice really requires – but to warrant this additional premise would amount to endorse JR). If this is right, Hindrik’s derivation of KR from KR-B is invalid.

Second, there is the *source objection*. Presumably, if we are entitled to criticise an unjustified assertion, it is not because the speaker has an unjustified belief. Recall GOOD COOK. If Bob’s unfounded claim (13) warrants criticism, it is not because Bob has formed an *unjustified belief*: as far as we’re concerned, Bob can believe whatever he wants. Instead, if we are entitled to criticise Bob it is because he unqualifiedly asserted that belief, despite being aware that he had no evidence in its support. Put differently, Bob is breaking a social norm (and thus liable to criticism) by making unjustified assertion, not by forming an unjustified belief. If this is right, the compound view mischaracterises what is objectionable about Bob’s claim, since

it traces the normative reason for our perception of an infraction (and our entitlement to criticise the statement) to Bob's unjustified belief.

3.5.5 The expansive view

Taken alone, the sincerity-rule doesn't capture what is objectionable about unjustified assertions. This, at least, if sincerity is understood to require speaker-centred sincerity, along BR's lines. But what about more "expansive" conceptions of insincerity, like unbound conceptions? Are there any that deem unjustified assertions insincere?

In GOOD COOK, Bob doesn't believe that what he claimed is false: he isn't insincere in the speaker-centred sense. However, Bob is aware that, by making a sincere statement, he could cause Rachel to acquire some other false belief. For example, Bob is aware that Rachel will likely infer that Bob has some evidence in support of his claim. In other words, Bob knows that his statement is *misleading*. To be sure, it isn't misleading in the narrow sense of being *intentionally designed* to deceive, but in the broad sense of having a *disposition to induce false beliefs*. Any account of insincerity that captures this "broad" notion of misleading will classify unjustified assertions as insincere.

In the spirit of hearer-centred views, the norm of sincerity could accordingly be understood as a duty to care about the effects that our statements have on our audiences – specifically, as a duty to ensure that one's assertions do not cause the audience to acquire false beliefs (cf. Eriksson 2011). The result would be an "expansive" belief-norm of assertion:

BR': One must: assert a proposition p only if one expects that one's assertion will not cause⁸⁶ the audience to accept⁸⁷ propositions that one believes to be false

This is a demanding conception of sincerity (since it is unbound), but not a particularly artificial one. An injunction like "do your best to avoid deceiving the audience" seems the kind of maxim that a speaker truly committed to sincerity should aim to follow. By this criterion, unjustified assertions like (13) are insincere. There is no need to defer this judgment to norms of belief, dodging the problems

⁸⁶ The "causing" should be understood to be determined by the speaker's making of an assertion, and not merely by their production of some sound (their locution). So, chanting to prove that one can sing, or reciting a line to illustrate one's acting ability, do not qualify as ways of violating BR'.

⁸⁷ To avoid the objections listed in §2.6.2, this can be expanded as "or continue to believe". This criterion is similar to Lackey's NMNA (cf. footnote XXX) – although Lackey presents it as a general communicative norm, not as a sincerity-rule.

faced by compound views. Bob is aware (as any competent speaker would be) that he might easily lead Rachel astray, causing her to believe some false proposition – for example, that he has stronger evidence for his claim than he actually possesses. Correspondingly, Bob should expect his statement to cause her to accept a false proposition.

What if Bob is oblivious to the effects that his statement might have on Rachel? One option is to insist that the statement is then impeccable, because it's meant to cause no harm. Another, more promising, is to suggest that, if sincerity is a duty to ensure that one's assertions don't cause the audience to acquire false beliefs, Bob's obliviousness by no means displays a truly sincere attitude. BR' then needs to be amended to rule out this option. A natural solution is to require that the speaker's expectation not to cause epistemic damage is actually reasonable.

BR'': One must: assert a proposition p only if one doesn't expect, with good reasons, that one's assertion will cause the audience to accept propositions that one believes to be false

The appearance of a requirement of "reasonableness" might raise the suspicion that BR'' borrows resources from justifications-rules (that define assertability in terms of reasonable belief)⁸⁸. This impression is mistaken. BR'' is a genuine (hearer-centred) sincerity-rule, because it only requires that the speaker aims not to cause a discrepancy between *their beliefs* and the *hearer's beliefs*. There is no requirement that the speaker's belief in their assertion (or in what they indirectly communicate with it) is reasonable. Only the assessment of the assertion's communicative effect on the speaker's audience must be reasonable.

A more pressing worry is that BR'-BR'' indirectly assume that assertions invite their audiences to believe that the speaker has evidence in their support – for example, Rachel can be expected to infer that Bob has evidence in support of (13). But if no norm like JR governs assertion, it's not obvious why hearers should be entitled to this inference. One could then suspect that BR'-BR'' implicitly assume JR, rather than representing an alternative to it.

However, pragmatic principles other than JR can explain why audiences are entitled to this conclusion. Consider the Maxim of Quantity, which dictates that contributions should not be less informative than required by the purpose of the exchange (cf. §3.2). The purpose of the exchange is to establish whether Jacques is a good cook. For this purpose, knowing that Bob's belief is based on a conjecture, instead of evidence of Jacques' ability, is certainly relevant information. Bob's (13) then violates the Maxim of Quantity unless he discloses this information (e.g. by

⁸⁸ A conception of sincerity that explicitly takes is in Eriksson (2011, 231–32).

qualifying his statement). This could entitle Rachel to her inference, with no need to appeal to JR⁸⁹.

If the foregoing is on the right track, BR'' can explain what is wrong about unjustified assertions without implicitly appealing to JR. For those who want to preserve the idea that sincerity, not justification, is the core norm of assertion, the expansive norm offers a great alternative to simple belief-rules like BR, even in their compound versions.

3.5.6 Only the facts? Truthfulness and assertion

The expansive view forbids unjustified assertions, and is therefore stronger than BR. Consequently, it fares better than BR, which was too lenient, delivering predictions similar to JR. But what if both BR'' and JR, despite being stricter than BR, still aren't strict enough? This is a standard accusation against justification-rules, which applies also to the expansive view. For both BR'' and JR, false assertions are perfectly permissible, as long as they are (e.g.) reasonably believed to be true.

Some philosophers contend that this take is problematic and too lenient. People share the *intuition* that (*ceteris paribus*) a false assertion is incorrect and improper. This intuition often translates into action: people typically criticise false assertions in virtue of their being false. Falsity itself constitutes a distinctive kind of wrongness for assertions. Any account failing to acknowledge that false assertions are incorrect and criticisable *in virtue of their being false* misses a fundamental linguistic datum about assertion. Factive accounts (like TR and KR) explain the wrongness and criticizability of falsity in terms of the violation of a factive norm. This explanation isn't open to non-factive accounts, like the BR and JR (cf. Williamson 2000:262).

If these considerations are correct, then assertion is subject to a norm of (objective) *truthfulness* on top of sincerity: to be epistemically permissible, a statement needs to be true, on top of being sincere. But is this right? Let's focus on non-factive rules that have survived previous scrutiny: JR and BR''. To be fair, also these views sometimes classify falsehoods as improper – for example, when they are disbelieved or unjustified. By contrast, when assertions are justifiedly believed but false (I'll call these UNLUCKY ASSERTIONS for simplicity), these views deem them appropriate. Is this verdict correct?

If UNLUCKY ASSERTIONS are impermissible in virtue of being false, factive accounts make the right call: objective truth is required for permissible assertion. By contrast, if UNLUCKY ASSERTIONS are permissible in virtue of being reasonably

⁸⁹ I presented one possible derivation, but more are possible – for example, by relying on the Relevance Principle (Wilson and Sperber 2002).

believed to be true, factive accounts are wrong: objective truth isn't required for permissible assertion. To test intuitions, let's consider an example (from Marsili and Wiegmann 2021):

COFFEE

Mallory manages an independent coffee shop. One of her customers is interested in the history and culture of coffee.

The customer asks Mallory whether the coffee is from Colombia. Mallory checks the coffee beans label, which says that the coffee is indeed from Colombia, so she replies:

(14) *The coffee is from Colombia*

However, Mallory doesn't know that the labels have been mixed up at the factory, and that the coffee beans are actually from Guatemala.

Mallory here says something false. But what else could have Mallory done? By giving the opposite answer ("The coffee is not from Colombia") Mallory would hit truth, purely by chance, and at the price of lying. By refusing to answer, she could surely avoid saying something false. But it would not only be rude: it presumably would violate some other communicative and epistemic norm (Goldberg 2020). On the other hand, responding with (14) seems totally appropriate, given what Mallory knows. Accordingly, many philosophers insist that UNLUCKY ASSERTIONS are permissible.

The intuition that Mallory would be making an appropriate assertion is overwhelmingly shared by non-philosophers: a strong majority (about 92%) of speakers disagrees with the statement "Mallory should not have said that the coffee is from Colombia". Generalising, empirical studies found that laypeople deem unlucky assertions appropriate and permissible (Kneer 2018; 2021; Reuter and Brössel 2019; Marsili and Wiegmann 2021).⁹⁰ Assuming that a good account of the norm of assertion "must face the linguistic data" (Douven 2006, 450; Kneer 2018), and make predictions that are generally consistent with the linguistic intuitions and appraisals of competent speakers, non-factive accounts are superior to factive ones⁹¹.

⁹⁰ Some initial findings by John Turri suggested that the norm of assertion is knowledge. However, every other researcher found incompatible results, and a strong case has now been made against the methodology employed in Turri's pioneering studies (Marsili and Wiegmann 2021; Graham and Pedersen 2024; Kneer and Marsili 2025).

⁹¹ Some factivist philosophers insist that these results are compatible with their view: if unlucky assertions are judged permissible, it is not because they comply with the norm, but rather because they violate it in an excusable way (comparable to the excusable promise-breaking discussed in §2.7.2). Apart from facing substantive theoretical objections (Douven 2006, 478–80; Lackey 2007; Gerken

To deny that the norm of assertion is factive is not to deny a link between assertion and truth (nor the importance of truthfulness for testimony). An essential link with truth can be maintained by arguing that truth is the *aim* (or *goal*) rather than the norm of assertion (Dummett 1973; Williams 2002; Marsili 2018b; 2021c). On this view, assertions are akin to other goal-directed activities: just like scoring a goal is the measure of success for a penalty shot in football (and failing to score a goal is not impermissible), truth is the measure of success for assertion (but failing to assert the truth is not *ipso facto* impermissible). An assertion “scores” and warrants positive evaluation when it represents reality accurately, and is defective (though not necessarily impermissible) when it fails to match the facts. Truth, in short, sets the standard for evaluating the success of assertions, not their permissibility.

This view explains why false assertions are intrinsically defective and criticisable, naturally complementing the normative picture predicated by justification-rules and sincerity-rules. If truth is the goal of assertion, it’s only natural that our social practice should forbid behaviour that is not conducive to asserting truths – such as making assertions that one doesn’t believe to be true, or that aren’t supported by adequate evidence. So understood, non-factive rules like BR” promote the achievement of assertion’s goal – their function is to maximise the proportion of true assertions. This, in turn, ensures that testimony is generally veridical.

3.6 The value of sincerity

This brings to a close our exploration of the normative dimension of sincerity. After briefly reviewing sincerity’s relation to other norms, its status as a virtue, and its role in supporting valuable social practices, I examined how a norm of sincerity might sustain assertion’s epistemic value. There’s a strong case in favour of the view that sincerity is required for permissible assertion. It seems, however, that epistemically permissible assertion requires more than sincerity, since it’s intuitively impermissible to assert *unjustified beliefs*.

Compound views tackle this limitation of sincerity-rules by supplementing sincerity with a separate norm, forbidding the acquisition of unjustified beliefs. However, since compound views face significant difficulties, I introduced an alternative: the *expansive view*. “Expansive” conceptions of sincerity forbid unjustified assertions indirectly, because unjustified assertions insincerely invite their audiences to believe

2011; Schechter 2017), this “excuse manoeuvre” faces empirical evidence that decisively speaks against it. Studies show that laypeople judge that unlucky assertions are plainly permissible, in patterns that deviate from the typical manifestations of excuse validation (Graham and Pedersen 2024; Kneer and Marsili 2025, 90–93).

that the speaker possesses adequate evidence. This view yields permissibility verdicts similar to the *justification-rule*, but reduces assertoric normativity entirely to expectations of sincerity. Paired with the idea that truth is the *aim* of assertion, both expansive sincerity-rules and justification-rules can acknowledge that false assertions are intrinsically defective, while avoiding the hosts of problems faced by factive rules like the knowledge-rule and the truth-rule. So conceived, non-factive views offer the best available explanations of the normativity of assertions.

The book's limited size means that many interesting questions about sincerity remain unexplored. Some open threads involve the application of our conceptual repertoire to more complex domains. For example, when multiple speakers (with multiple beliefs) jointly make an assertion as a group, what determines whether they are sincere?⁹² Other questions concern the study of sincerity in other areas of philosophy (like moral or political philosophy), in other disciplines (social and experimental psychology, politics, and the law) and in other philosophical and cultural traditions (e.g. Rogacz 2022). The ideas and arguments developed here are therefore not intended as the final word, but as tools and resources for better understanding sincerity, both in everyday life and in scholarly pursuits.

⁹² For discussion, Lackey (2020) and Marsili (2023).

4. References

- Adler, Jonathan E. 2006. 'Epistemological Problems of Testimony'. In *Stanford Encyclopedia of Philosophy*, Winter 2017. <https://plato.stanford.edu/archives/win2017/entries/testimony-episprob/>.
- Aldrich, Virgil C. 1966. 'Telling, Acknowledging and Asserting'. *Analysis* 27 (2): 53–56.
- Alfano, Mark, and Neil Levy. 2020. 'Knowledge From Vice: Deeply Social Epistemology.' *Mind* 129 (515): 887–915.
- Alston, William P. 2000. *Illocutionary Acts and Sentence Meaning*. Ithaca: Cornell University Press.
- Aquinas. ST. *Summa Theologica*. Lulu.com.
- Aristotle. NE. *Nicomachean Ethics*. Hackett Publishing.
- Augustine. DM. *De Mendacio*. CreateSpace Independent Publishing Platform.
- Austin, John Langshaw. 1975. *How To Do Things With Words*. 2nd ed. Oxford: Clarendon Press.
- Bach, Kent. 2008. 'Applying Pragmatics to Epistemology'. *Philosophical Issues*, no. 18, 68–88.
- Bach, Kent, and Robert M. Harnish. 1979. *Linguistic Communication and Speech Acts*. Cambridge: MIT Press.
- Benton, Matthew Aaron. 2018. 'Lying, Accuracy and Credence'. *Analysis* 78 (2): 195–98.
- Bok, Sissela. 1978. *Lying*. Random House.
- Carson, Thomas L. 2006. 'The Definition of Lying'. *Noûs* 40 (2): 284–306.
- . 2010. *Lying and Deception*. Oxford: Oxford University Press.
- . 2018. 'Lying and Ethics'. In *The Oxford Handbook of Lying*, edited by Jörg Meibauer, 469–82.
- Carson, Thomas L., Richard E Wokutch, and Kent F Murrmann. 1982. 'Bluffing in Labor Negotiations: Issues Legal and Ethical'. *Journal of Business Ethics* 1 (1): 13–22.
- Chan, Timothy, and Guy Kahane. 2011. 'The Trouble with Being Sincere'. *Canadian Journal of Philosophy* 41 (2): 1–13. <http://www.ncbi.nlm.nih.gov/pmc/articles/pmc3272424/>.
- Chisholm, Roderick M., and Thomas D Feehan. 1977. 'The Intent to Deceive'. *Journal of Philosophy* 74 (3): 143–59.
- Clem, Stewart. 2023. *Lying and Truthfulness: A Thomistic Perspective*. Cambridge: Cambridge University Press.
- Coady, C A J. 1992. *Testimony*. Oxford: Oxford University Press.
- Cohen, G A. 2002. 'Deeper into Bullshit' 1.

- Cullison, Andrew. 2010. 'On the Nature of Testimony'. *Episteme*, 114–27.
- Davis, Wayne. 1999. 'Communicating, Telling and Informing'. *Philosophical Inquiry* 21 (1): 21–43.
- . 2003. *Meaning, Expression and Thought*. Cambridge: Cambridge University Press.
- . 2010. 'Implicature'. In *Stanford Encyclopedia of Philosophy*, 1–37.
- Douven, Igor. 2006. 'Assertion, Knowledge, and Rational Credibility'. *The Philosophical Review* 115 (4): 449–85.
- Dummett, Michael. 1973. 'Assertion'. In *Frege: Philosophy of Language*, edited by Duckworth.
- Erickson, Thomas D., and Mark E. Mattson. 1981. 'From Words to Meaning: A Semantic Illusion'. *Journal of Verbal Learning and Verbal Behavior* 20 (5): 540–51.
- Eriksson, John. 2011. 'Straight Talk: Conceptions of Sincerity in Speech'. *Philosophical Studies* 153 (2): 213–34.
- Falkenberg, Gabriel. 1988. 'Insincerity and Disloyalty'. *Argumentation* 2 (1): 89–97.
- Fallis, Don. 2010. 'Lying and Deception'. *Philosophers' Imprint* 10 (11).
- . 2012. 'Lying as a Violation of Grice's First Maxim of Quality'. *Dialectica* 66 (4): 563–81.
- . 2018. 'Lying and Omissions'. In *The Oxford Handbook of Lying*, 183–92. Oxford University Press.
- Faulkner, Paul. 2013. 'Lying and Deceit'. In *The International Encyclopedia of Ethics*. <http://onlinelibrary.wiley.com/doi/10.1002/9781444367072.wbiec482/full>.
- Foucault, M. 2011. *The Courage of Truth*. Translated by Graham Burchell. 2011th edition. Palgrave Macmillan.
- Frankfurt, Harry. 1986. *On Bullshit*. Princeton, NJ: Princeton University Press.
- Frege, Gottlob. 1956. 'The Thought: A Logical Inquiry'. *Mind*. Oxford University Press/Mind Association.
- García-Carpintero, Manuel. 2004. 'Assertion and the Semantics of Force-Markers'. In *The Semantics/Pragmatics Distinction*, edited by Claudia Bianchi, 133–166.
- . 2022. 'How to Understand Rule-Constituted Kinds'. *Review of Philosophy and Psychology* 13 (1): 7–27.
- Gerken, Mikkel. 2011. 'Warrant and Action'. *Synthese* 178 (3): 529–47.
- Goldberg, Sanford C. 2011. 'Putting the Norm of Assertion to Work: The Case of Testimony'. In *Assertion: New Philosophical Essays*, edited by Jessica Brown and Herman Cappelen, 1–30. Oxford: Oxford University Press.
- . 2015. *Assertion: On the Philosophical Significance of Assertoric Speech*. Oxford University Press.
- Goldberg, Sanford C. 2020. *Conversational Pressure: Normativity in Speech Exchanges*. Oxford University Press.
- González De Prado, Javier. 2023. 'No Norm for (off the Record) Implicatures'. *Inquiry*, August, 1–21.
- Graham, Peter J. 2020. 'Assertions, Handicaps, and Social Norms'. *Episteme* 17 (3): 349–63.

- Graham, Peter J., and Nikolaj J. L. L. Pedersen. 2024. 'Knowledge Is Not Our Norm of Assertion'. In *Contemporary Debates in Epistemology (3rd Edition)*, edited by Blake Roeber, John Turri, Matthias Steup, and Ernest Sosa. New York: Routledge.
- Green, Adam. 2017. 'An Epistemic Norm for Implicature'. *Journal of Philosophy* 114 (7): 381–91.
- Green, Mitchell. 2007a. *Self-Expression*. Oxford: Oxford University Press.
- . 2007b. 'Speech Acts'. In *Stanford Encyclopedia of Philosophy*, 1–19.
- Green, Stuart P. 2018. 'Lying and the Law'. In *The Oxford Handbook of Lying*, edited by Jörg Meibauer, 0. Oxford University Press.
- Grice, H. P. 1957. 'Meaning'. *The Philosophical Review* 66 (3): 377–88.
- . 1989. *Studies in the Way of Words*. Cambridge, MA: Harvard University Press.
- Griffiths, Paul J. 2004. *Lying: An Augustinian Theology of Duplicity*. Brazos Press.
- Güngör, Hüseyin. 2024. 'Non-Literal Lies Are Not Exculpatory'. *The Philosophical Quarterly*, July, pqa078.
- Hindriks, Frank. 2007. 'The Status of the Knowledge Account of Assertion'. *Linguistics and Philosophy* 30 (3): 393–406.
- Hoek, Daniel. 2018. 'Conversational Exculpature'. *Philosophical Review* 127 (2): 151–96.
- Hughes, Nick. 2017. 'No Excuses: Against the Knowledge Norm of Belief'. *Thought: A Journal of Philosophy*, July.
- Isenberg, Arnold. 1964. 'Deontology and the Ethics of Lying'. *Philosophy and Phenomenological Research* 24 (4): 463–80.
- Kant, Immanuel. 1797. 'On a Supposed Right to Lie Because of Philanthropic Concerns', no. 1.
- . LE. *Lectures on Ethics*. Edited by Peter Heath and J. B. Schneewind. Translated by Peter Heath. The Cambridge Edition of the Works of Immanuel Kant. Cambridge: Cambridge University Press.
- Kneer, Markus. 2018. 'The Norm of Assertion: Empirical Data'. *Cognition* 177 (July 2017): 165–71.
- . 2021. 'Norms of Assertion in the United States, Germany, and Japan'. *PNAS* 118 (37): 3.
- Kneer, Markus, and Neri Marsili. 2025. 'The Truth about Assertion and Retraction: A Review of the Empirical Literature'. In *Lying, Fake News, and Bullshit*, edited by Alex Wiegmann. Bloomsbury.
- Korsgaard, Christine M. 1986. 'The Right to Lie: Kant on Dealing with Evil'. *Philosophy Public Affairs* 15 (4): 325–49.
- Krauss, Sam Fox. 2017. 'Lying, Risk and Accuracy'. *Analysis* 73:651–59.
- Krstić, Vladimir. 2019. 'Can You Lie Without Intending to Deceive?' *Pacific Philosophical Quarterly* 100 (2): 642–60.
- . 2023. 'Lying: Revisiting the "Intending to Deceive" Condition'. *Analysis*, April, anac099.
- Lackey, Jennifer. 2006. 'The Nature of Testimony'. *Pacific Philosophical Quarterly* 87 (2): 177–97.
- . 2007. 'Norms of Assertion'. *Noûs* 41 (4): 594–626.

- . 2020. *The Epistemology of Groups*. Oxford University Press.
- Leonard, Nick. 2023. ‘Epistemological Problems of Testimony’. In *The Stanford Encyclopedia of Philosophy*, edited by Edward N. Zalta and Uri Nodelman, Spring 2023. Metaphysics Research Lab, Stanford University. <https://plato.stanford.edu/archives/spr2023/entries/testimony-episprob/>.
- Lewis, David K. 1975. ‘Languages and Language’. In *Minnesota Studies in the Philosophy of Science*, edited by Keith Gunderson, 3–35. University of Minnesota Press.
- Mahon, James Edwin. 2006. ‘Kant and the Perfect Duty to Others Not to Lie’. *British Journal for the History of Philosophy* 14 (4): 653–85.
- . 2015. ‘The Definition of Lying and Deception’. In *Stanford Encyclopedia of Philosophy*.
- Mandelkern, Matthew, and Kevin Dorst. 2022. ‘Assertion Is Weak’. *Philosophers’ Imprint*.
- Marsili, Neri. 2014. ‘Lying as a Scalar Phenomenon’. In *Certainty-Uncertainty – and the Attitudinal Space in Between*, edited by Sibilla Cantarini, Werner Abraham, and Elisabeth Leiss, 153–73. Amsterdam: John Benjamins Publishing Company.
- . 2016. ‘Lying by Promising’. *International Review of Pragmatics* 8 (2): 271–313.
- . 2017. ‘You Don’t Say ! Lying, Asserting and Insincerity’. PhD Dissertation, University of Sheffield. <https://etheses.whiterose.ac.uk/19068/>.
- . 2018a. ‘Lying and Certainty’. In *The Oxford Handbook of Lying*, edited by Jörg Meibauer, 169–82. Oxford: Oxford University Press.
- . 2018b. ‘Truth and Assertion: Rules versus Aims’. *Analysis* 78 (4): 638–48.
- . 2021a. ‘Eliot Michaelson and Andreas Stokke (Eds.), Lying: Language, Knowledge, Ethics, and Politics (Oxford: Oxford University Press, 2018), Pp. 320.’ *Utilitas* 33 (4): 502–5.
- . 2021b. ‘Lying, Speech Acts, and Commitment’. *Synthese* 199:3245–69.
- . 2021c. ‘Truth: The Rule or the Aim of Assertion?’ *Episteme*, September, 1–7.
- . 2022. ‘Immoral Lies and Partial Beliefs’. *Inquiry* 65 (1): 117–27.
- . 2023. ‘Group Assertions and Group Lies’. *Topoi* 42 (2): 369–84.
- . 2025. ‘How Online Misinformation Works: A Costly Signalling Perspective, to Appear In P. (Ed.)’: In *Misinformation and Other Epistemic Pathologies.*, edited by M. Popa-Wyatt. Cambridge University Press.
- Marsili, Neri, and Mitchell Green. 2021. ‘Assertion: A (Partly) Social Speech Act’. *Journal of Pragmatics* 181 (August):17–28.
- Marsili, Neri, and Guido Löhr. 2022. ‘Saying, Commitment, and the Lying-Misleading Distinction’. *The Journal of Philosophy* 119 (12): 687–98.
- Marsili, Neri, and Alex Wiegmann. 2021. ‘Should I Say That? An Experimental Investigation of the Norm of Assertion.’ *Cognition* 212.
- Mazzarella, Diana. 2023. ‘“I Didn’t Mean to Suggest Anything like That!”: Deniability and Context Reconstruction’. *Mind & Language* 38 (1): 218–36.
- McGlynn, Aidan. 2014. *Knowledge First?* Basingstoke: Palgrave Macmillan.
- Meibauer, Jörg. 2005. ‘Lying and Falsely Implicating’. *Journal of Pragmatics* 37 (9): 1373–99.

- . 2014. *Lying at the Semantics-Pragmatics Interface*. *Lying at the Semantics-Pragmatics Interface*. Berlin, Boston: De Gruyter.
- Mellor, D. H. 1977. 'Conscious Belief'. *Proceedings of the Aristotelian Society* 78:87–101. <https://www.jstor.org/stable/4544919>.
- Milić, Ivan. 2017. 'Against Selfless Assertions'. *Philosophical Studies* 174 (9): 2277–95.
- Miller, Christian B. 2021. *Honesty: The Philosophy and Psychology of a Neglected Virtue*. New York: Oxford University Press.
- Millikan, Ruth Garrett. 2005. *Language: A Biological Model*. Oxford: Clarendon Press.
- Moran, Richard. 2005. 'Problems of Sincerity'. *Proceedings of the Aristotelian Society* 105 (1): 325–45.
- Owens, David. 2006. 'Testimony and Assertion'. *Philosophical Studies* 130 (1): 105–29.
- Pagin, Peter, and Neri Marsili. 2021. 'Assertion'. In *Stanford Encyclopedia of Philosophy*, Winter 2021 edition. <https://plato.stanford.edu/archives/win2021/entries/assertion/>.
- Pennycook, Gordon, James Allan Cheyne, Nathaniel Barr, Derek J Koehler, and Jonathan A Fugelsang. 2015. 'On the Reception and Detection of Pseudo-Profound Bullshit'. *Judgment and Decision Making* 10 (6): 549–63.
- Pepp, Jessica. 2018. 'Truth Serum, Liar Serum, and Some Problems about Saying What You Think Is False'. In *Lying: Language, Knowledge, Ethics, Politics*. Oxford University Press.
- . n.d. 'The Size of a Lie: From Truthlikeness to Sincerity'. *Inquiry* 0 (0): 1–24. Accessed 19 July 2024.
- Pinker, Steven, Martin A. Nowak, and James J. Lee. 2008. 'The Logic of Indirect Speech'. *Proceedings of the National Academy of Sciences of the United States of America* 105 (3): 833–38.
- Queloz, Matthieu. 2021. *The Practical Origins of Ideas: Genealogy as Conceptual Reverse-Engineering*. 1st ed. Oxford University Press Oxford.
- Reuter, Kevin, and Peter Brössel. 2019. 'No Knowledge Required'. *Episteme* 16 (3): 303–21.
- Ridge, Michael. 2006. 'Sincerity and Expressivism'. *Philosophical Studies* 131 (2): 487–510.
- Rogacz, Dawid. 2022. 'Sincerity (Cheng) as a Civic and Political Virtue in Classical Confucian Philosophy'. *Philosophy Compass* 17 (6): e12833.
- Rousseau, Jean-Jacques. RSW. *The Reveries of the Solitary Walker*. Hackett Publishing.
- Russell, Bertrand. 1946. *History of Western Philosophy*. Routledge.
- Rutschmann, Ronja, and Alex Wiegmann. 2017. 'No Need for an Intention to Deceive? Challenging the Traditional Definition of Lying'. *Philosophical Psychology*.
- Saul, Jennifer M. 2012. *Lying, Misleading, and What Is Said: An Exploration in Philosophy of Language and Ethics*. Oxford University Press.
- Saul, Jennifer M. 2024. *Dogwhistles and Figleaves: How Manipulative Language Spreads Racism and Falsehood*. Oxford, New York: Oxford University Press.
- Schechter, Joshua. 2017. 'No Need for Excuses Against Knowledge-First Epistemology and the Knowledge Norm of Assertion'. In *Knowledge First:*

- Approaches in Epistemology and Mind*, edited by J. Adam Carter, Emma C. Gordon, and Benjamin W. Jarvis, 1–295. Oxford University Press.
- Schiffer, Stephen R. 1972. *Meaning*. Oxford; Clarendon Press.
- Searcy, William, and Stephen Nowicki. 2005. *The Evolution of Animal Communication: Reliability and Deception in Signaling Systems*. Princeton University Press.
- Searle, John R. 1969. *Speech Acts: An Essay in the Philosophy of Language*. Cambridge University Press.
- . 1979. *Expression and Meaning*. Cambridge: Cambridge University Press.
- Sedgwick, Sally. 1991. ‘On Lying and the Role of Content in Kant’s Ethics’. *Kant-Studien* 82 (1): 42–62.
- Siebel, Mark. 2003. ‘Illocutionary Acts and Attitude Expression’. *Linguistics and Philosophy* 26 (3): 351.
- . 2020. ‘The Belief View of Assertion’. *The Oxford Handbook of Assertion*, 97–118.
- Siegler, Frederick A. 1966. ‘Lying’. *American Philosophical Quarterly* 3 (2): 128–36.
- Skyrms, Brian. 2010. *Signals: Evolution, Learning, & Information*. Oxford; New York: Oxford University Press.
- Smith, John Maynard, and David Harper. 2003. *Animal Signals*. Oxford Series in Ecology and Evolution. Oxford, New York: Oxford University Press.
- Sneddon, Andrew. 2020. ‘Alternative Motivation and Lies’. *Analysis*, September.
- Sorensen, Roy. 2007. ‘Bald-Faced Lies! Lying without the Intent to Deceive’. *Pacific Philosophical Quarterly* 88:251–64.
<http://onlinelibrary.wiley.com/doi/10.1111/j.1468-0114.2007.00290.x/full>.
- . 2010. ‘Knowledge-Lies’. *Analysis* 70 (4): 608–15.
- . 2011. ‘What Lies behind Misspeaking’. *American Philosophical Quarterly* 48 (4): 399–410.
- . 2018. ‘Lying to Mindless Machines’. In *Lying: Language, Knowledge, Ethics, Politics*, edited by Eliot Michaelson and Andreas Stokke, 1–23. Oxford University Press.
- . 2023. ‘Vagueness’. In *The Stanford Encyclopedia of Philosophy*, edited by Edward N. Zalta and Uri Nodelman, Winter 2023. Metaphysics Research Lab, Stanford University.
<https://plato.stanford.edu/archives/win2023/entries/vagueness/>.
- Sperber, Dan, Fabrice Clement, Christophe Heintz, Olivier Mascaro, Hugo Mercier, Gloria Origgi, and Deirdre Wilson. 2010. ‘Epistemic Vigilance’. *Mind and Language* 25 (4): 359–93.
- Stenius, Erik. 1967. ‘Mood and Language-Game’. *Synthese* 17 (3): 254–74.
<https://www.jstor.org/stable/20114558>.
- Stokke, Andreas. 2018. *Lying and Insincerity*. Oxford: Oxford University Press.
- Strawson, P. F. 1964. ‘Intention and Convention in Speech Acts’. *The Philosophical Review* 73 (4): 439–60.
- Strudler, Alan. 2009. ‘The Distinctive Wrong in Lying’. *Ethical Theory and Moral Practice* 13 (2): 171–79.
- Swift, Jonathan. 1710. *The Art Of Political Lying*. <http://archive.org/details/swift-the-art-of-political-lying-1710>.

- Tagliapietra, Andrea. 2003. *La virtù crudele. Filosofia e storia della sincerità*. Einaudi.
- Trilling, Lionel. 2009. *Sincerity and Authenticity*. Harvard University Press.
- Turri, John. 2014. 'Selfless Assertions: Some Empirical Evidence'. *Synthese*, no. October 2014, 1–23.
- Varga, Somogy, and Charles Guignon. 2023. 'Authenticity'. In *The Stanford Encyclopedia of Philosophy*, edited by Edward N. Zalta and Uri Nodelman, Summer 2023. Metaphysics Research Lab, Stanford University. <https://plato.stanford.edu/archives/sum2023/entries/authenticity/>.
- Viebahn, Emanuel. 2017. 'Non-Literal Lies'. *Erkenntnis*.
- . 2021. 'The Lying/Misleading Distinction: A Commitment-Based Approach'. *Journal of Philosophy* CXVIII (6).
- Weiner, Matthew. 2005. 'Must We Know What We Say?' *The Philosophical Review* 114 (2): 227–51.
- Whiting, Daniel. 2012. 'Stick to the Facts: On the Norms of Assertion'. *Erkenntnis* 78 (4): 847–67.
- Williams, Bernard Arthur Owen. 2002. *Truth and Truthfulness An Essay in Genealogy*. Princeton: Princeton University Press.
- Wilson, Deirdre. 1995. 'Is There a Maxim of Truthfulness?'. *UCL Working Papers in Linguistics*, no. 7, 197–212.
- Wilson, Deirdre, and Dan Sperber. 2002. 'Truthfulness and Relevance'. *Mind* 25:1–41.