

# What Is Knowledge?

*You looked at the clock. You were right. Did you know?*



● STOPPED 12 H AGO – BUT  
RIGHT, FOR THIS ONE MINUTE

It is 9:12 in the morning and you are late. You glance up at the great station clock as you rush past, read **9:12**, and think: *fine – three minutes to spare*. You are right. It really is 9:12. And yet the clock you trusted died at 9:12 exactly twelve hours ago, somewhere in the small hours, and has hung there frozen ever since. You consulted a broken instrument at the single instant in the day it happened to be correct.

Your belief was **true**. It rested on a perfectly sensible **reason** – clocks tell the time, and you have trusted a thousand of them without incident. You **believed** it sincerely. So: did you *know* it was 9:12? Asked carefully, almost everyone says no. Something is missing. Saying exactly what has consumed philosophers for sixty years – and, as we'll see, the better part of a thousand.

This is the first descent, so there is nothing behind us yet – the log is blank. Instead we plant seeds. The machinery introduced today (belief as something that comes in *degrees*; updating on evidence; minds as inference engines) is the epistemic toolkit the entire course will lean on. Watch for it to resurface on **Day 2** (how science decides what counts at all), **Day 4** (probability as the logic of partial belief), **Day 7** (information), **Day 119** (the predictive brain), and **Day 149** (when famous results evaporate). The five threads we'll trace across all 180 days – *information, energy, evolution, emergence, computation* – all have a quiet first appearance right here.

## — THE MODEL

### The three-legged stool

For roughly twenty-three centuries, Western philosophy carried around a tidy answer to "what is knowledge?" To *know* that something is the case, you needed three things at once:

**(1) you believe it** – you can't know what you don't even hold to be true. **(2) it's true** – you can't *know* a falsehood; people who said "I knew the Earth was flat" merely *believed* it, confidently and wrongly. **(3) you're justified** – you have good reason, because a lucky guess that lands isn't knowledge either. The gambler who "just had a feeling" the long-shot would win, and won, did not *know* it would.

Knowledge, on this view, is *justified true belief* – JTB, a three-legged stool. Kick away any leg and it topples. The picture is usually traced to Plato, who in the *Theaetetus* floats the idea that knowledge is "true judgement with an account." There's a delicious irony here, much enjoyed by historians: in that very dialogue Socrates then dismantles the definition, so Plato arguably never endorsed the thing named after him. As one scholar put it, it is almost as if a distinguished critic created a tradition in the very act of destroying it.

Still, the rough consensus held. The stool seemed stable. And then a 35-year-old philosopher who, the story goes, hadn't published much and rather needed to, wrote three pages.

## — THE GRENADE

## Gettier's three pages

In 1963, Edmund Gettier published a paper in the journal *Analysis* with the cheekily plain title "*Is Justified True Belief Knowledge?*". It runs barely three pages. It has since been cited in **thousands** of scholarly works and spawned entire subfields. Few documents in modern philosophy have done more damage per word.

Gettier's move was devastatingly simple. He built little stories in which all three legs of the stool are firmly in place – belief, truth, justification – and yet you'd never say the person *knows*. Here is his first case, lightly modernized:

*Smith and Jones both apply for a job. The boss tells Smith, "Jones will get it." Smith has also, idly, counted the coins in Jones's pocket: ten. So Smith forms a justified belief: the person who gets the job has ten coins in their pocket.*

Now the twist. The boss was wrong (or changed her mind): **Smith** gets the job, not Jones. And – entirely unknown to Smith – Smith happens to have **ten coins** in his own pocket too. Look at his belief, "the person who gets the job has ten coins": it's **true** (the winner, Smith, does have ten coins), it's **justified** (excellent evidence – the boss's word, a literal coin count), and it's sincerely **believed**. JTB, all three legs. Yet Smith plainly doesn't *know* it. He was tracking *Jones* and arrived at the right answer about the wrong man.

That is the anatomy of a *Gettier case*: your justification runs *through a falsehood* ("Jones will get the job"), and the belief is rescued into truth by an unrelated *coincidence* ("Smith also has ten coins"). The reason and the truth never actually touch. The stopped clock is the same skeleton in cleaner clothes: your reason (the clock) is broken, and the truth (it's 9:12) arrives by luck.

A TWIST OLDER THAN ITS NAME

Gettier wasn't first. Bertrand Russell had the stopped-clock case in *Human Knowledge: Its Scope and Limits* (1948). Go back further and the problem is downright ancient: around **770 CE** the Buddhist logician **Dharmottara** described a traveler who sees what looks like smoke over a hill, infers fire, and is right that there's fire – except the "smoke" was a swarm of insects. Same skeleton, twelve centuries early. In 14th-century India, **Gaṅgeśa** built a whole causal theory of knowing to handle such cases. The "Gettier problem" is one of philosophy's great instances of *convergent discovery* – the kind of thing minds keep tripping over independently, which is itself a hint that something real is there.

### The Gettier Machine

CASE	BELIEF	TRUTH	JUSTIFICATION	LUCK	VERDICT
Plain knowing	Yes	Yes	Yes	No	Knowledge on the classic view
Stopped clock	Yes	Yes	Yes	Yes	Not knowledge: truth arrives by coincidence
Lucky guess	Yes	Yes	No	Yes	Not knowledge: no justification
Confident error	Yes	No	Yes	No	Not knowledge: the claim is false

— THE PATCH WARS

## The hunt for the fourth leg

The obvious response to Gettier was: add a fourth condition that screens out the luck. For decades, epistemologists tried – and each tidy fix met a nastier counterexample. It became a minor blood sport.

**No false lemmas.** First idea: knowledge can't be reasoned *through* a falsehood. Smith's belief leaned on "Jones will get the job," which was false; ban that and you're safe. Clean – until Alvin Goldman's **fake-barn country** (1976). You're driving through a region where, as a prank, every "barn" is a flat movie-set façade – except one. You happen to glance at the single real barn and think "a barn." Your belief is true, justified, and rests on *no* false premise. Yet you don't know it's a barn: you could so easily have been fooled by a façade a hundred meters either way.

**Track the truth.** So maybe knowledge is about how your belief behaves across *nearby possibilities*. Robert Nozick (1981) proposed *sensitivity*: you know *p* only if, *were p false, you wouldn't believe it*. Elegant – but it produces strange verdicts in edge cases. Ernest Sosa (1999) flipped it into *safety*: in all the nearby ways things could have gone, you wouldn't have been wrong. The stopped clock fails safety hard (a minute either side and you're mistaken); a working clock passes. Fake-barn-you fails safety too.

Then Linda Zagzebski (1994) delivered the gut-punch with a kind of **recipe** for defeating *any* such fix. Take a belief that's justified but could still be false (which justification, being fallible, always allows). Arrange for the justification to misfire so the belief is false – then arrange, by luck, for it to be true after all. As long as your fourth condition stops short of demanding that the justification *guarantee* the truth, luck can always wedge back in. The patch wars may be structurally unwinnable.

## Two ways to stop fighting

**Declare knowledge a primitive.** Timothy Williamson, in *Knowledge and Its Limits* (2000), made a radical move: stop trying to build knowledge out of simpler parts. Maybe it has no analysis. On his *knowledge-first* view, knowing is a basic mental state – the most general *factive* one – and we should explain belief, evidence, and justification *in terms of knowledge*, not the other way around. You can't define *hydrogen* or *John F. Kennedy* into simpler concepts; perhaps knowledge is bedrock too. Sixty years of failed definitions start to look less like a puzzle and more like a clue.

**Make it about competence.** The other escape is *virtue epistemology* (Sosa again). Knowledge is *apt* belief – a belief that is true *because of* the knower's skill, not by accident. Picture an archer. A bullseye is a good shot only if the arrow hit center *because* the archer aimed well – not because a gust blew a bad shot onto the target. The Gettiered believer is exactly that archer: the wind knocked the arrow off course, then a second gust knocked it back onto the bull. Accurate, yes. Skillful, no. *Apt*, no. That, says Sosa, is why luck-based hits aren't knowledge.

## — THE DEBATE

## What makes a belief justified at all?

Step back from "is it knowledge?" to the humbler leg: what makes a belief *justified* in the first place? Push on any justification and you fall into a regress. It's 9:12 because the clock says so. Trust the clock because clocks are reliable. Believe *that* because... and now you're sliding. The ancient skeptics mapped the trap precisely. Every chain of justification, they argued, ends in one of three uncomfortable places – the *Agrippan trilemma*: it goes on **forever**, or it loops back in a **circle**, or it stops at some **arbitrary** point you simply declare.

Three modern schools each pick which horn to grab – and a fourth changes the subject entirely.

### DIAGRAM · THE REGRESS PROBLEM

## Agrippa's Trilemma — three bad endings, four escapes

Why is your belief justified? Every honest answer to "...and why *that*?" eventually hits one of three walls.

**Reason chain:** belief: "it's 9:12" -> because "the clock" -> because "...and why that?"

1. **Infinite regress:** every reason needs another reason forever.
2. **Circle:** the chain loops back to something it already used.
3. **Arbitrary halt:** the chain simply stops at a basic commitment.

**Foundationalism** – bites the third bullet: some beliefs are *basic* and need no further support (raw experience, simple logic). The chain stops, but not arbitrarily.

**Infinetism** – the brave minority: accepts that justification is an endless chain of reasons, never bottoming out.

**Coherentism** – embraces the circle, but makes it virtuous: no belief stands alone; a belief is justified by how well it hangs together with the whole web. (A first taste of *systems thinking*, Day 9.)

**Reliabilism** – changes the question. A belief is justified if it was *produced by a reliable process* – good vision, sound memory – whether or not you can recite a defense. This is *externalism*: justification can be a fact about your wiring, not a story in your head.

That internal/external split matters more than it looks. The **internalist** says justification must be something you can access by reflection – reasons available "from the inside." The **externalist** (reliabilism's home) says what matters is that your belief was, in fact, produced in a truth-conducive way, accessible or not. Hold that tension in mind: it is exactly where the old armchair questions collide with the new science of how brains actually form beliefs.

---

— THE FRONTIER · 2026

## Three live edges — and the hype filter

Every day in this course ends at the research frontier, with each claim tagged for how much weight it can bear. Knowledge sits at a fascinating junction right now: philosophers, psychologists, and neuroscientists are all circling the same questions from different sides.

---

Edge 01 [SUPERSEDED] [ESTABLISHED]

### Are "knowledge" intuitions universal — or just Western?

When the discipline runs on "asked carefully, almost everyone says no," a natural worry is: *which* everyone? In 2001, the founding study of *experimental philosophy* – Weinberg, Nichols & Stich – reported that the Gettier intuition varies by culture, with East-Asian participants supposedly more willing to grant the lucky believer "knowledge." If true, it was a bombshell: philosophy's whole method of consulting intuitions looked parochial.

The bombshell did not survive contact with replication. In "**Gettier Across Cultures**" (*Notis*, 2017), Machery, Stich, Rose and colleagues tested Brazil, India, Japan, and the United States with cases taken near-verbatim – and found the *opposite*: in **every** group, people robustly refused to call the Gettiered belief knowledge. A separate replication (Kim & Yuan) failed to reproduce the original cross-cultural gap even with a far larger East-Asian sample. The current best reading is that there may be a **universal core "folk epistemology"** that recoils from luck-based knowing. The deeper lesson is one we'll meet at industrial scale on **Day 149**: the splashiest finding is often the one careful re-testing quietly walks back.

---

Edge 02 [ESTABLISHED] [CONTESTED]

## Belief by the dial, not the switch: Bayesian epistemology

Maybe the all-or-nothing picture of belief was the wrong starting point. *Bayesian epistemology* says your real epistemic states are *credences* – degrees of confidence on a scale from 0 to 1. Rationality then needs just two rules: your credences must obey the laws of probability (*coherence*), and you must revise them by *conditionalization* as evidence comes in.

Why obey? The **Dutch book theorem** (Ramsey, 1926; de Finetti, 1937) supplies a startlingly concrete answer: if your credences break the probability laws, a clever bookmaker can offer you a set of bets you'll each accept as fair, but which together guarantee you lose money *no matter what happens*. Incoherent confidence isn't merely untidy – it's exploitable. The dial below lets you feel the trap close. What's still *contested* is whether graded credence *replaces* ordinary yes/no belief or merely sits beside it. (The lottery paradox bites here: you're 99.9% sure your ticket loses – but do you flat-out *believe* it loses?) We pick this thread up properly on **Day 4**.

### The Credence Dial and the Dutch Book

If your credence in  $S$  and your credence in  $not-S$  sum to 1.00, the pair is coherent. If they sum above 1.00, you will overpay for bets where exactly one can win. If they sum below 1.00, a bookie can reverse the bets and still guarantee a profit.

CREDENCE IN S	CREDENCE IN NOT-S	SUM	RESULT
0.50	0.50	<b>1.00</b>	Coherent
0.70	0.60	<b>1.30</b>	Guaranteed 0.30 loss if you buy both \$1 bets
0.30	0.40	<b>0.70</b>	Guaranteed 0.30 loss if the bookie buys both bets from you

Edge 03 [PROMISING] [CONTESTED]

## Where do beliefs come from? The brain as a prediction machine

Philosophy asks what justifies a belief; neuroscience now asks how a lump of tissue forms one. A fast-growing program answers: the brain is not a passive sponge soaking up the world – it is a relentless *prediction machine*. On the *predictive-processing* view (Andy Clark, *Behavioral and Brain Sciences*, 2013; Jakob Hohwy, 2013), the brain constantly generates a model of its surroundings, predicts the sensory signals it expects, and forwards only the *prediction errors* – the surprises – up the hierarchy. Perception becomes the brain's best running guess, reined in by error; in Anil Seth's memorable phrase, a "controlled hallucination." Belief-updating starts to look like **Bayesian inference rendered in neurons** – the so-called "Bayesian brain," tying Edge 02 to wetware.

Karl Friston pushes the idea to its limit with the *Free Energy Principle* (*Nature Reviews Neuroscience*, 2010): living systems persist precisely by minimizing a quantity – "free energy," an information-theoretic cousin of *surprise* – that knits perception, action, and even biological self-organization into one framework. The honest labels matter here. Predictive coding genuinely explains real perceptual phenomena and is a serious, productive research program – **promising**. But the *grand* Free Energy Principle, as a single law for all of mind and life, is widely criticized as so general it is hard to *falsify* – closer to a framework than a tested theory, and so **contested**. We'll return to it for perception (**Day 119**) and consciousness (**Days 123–126**) – and notice already how its "free energy" rhymes with the thermodynamics we'll meet on **Days 33 and 83–85**. *Information, energy, computation, emergence* – four of our five threads, braided into one neuron's quiet arithmetic.

### — OPEN QUESTIONS

## What's genuinely unsettled

Sixty years on, the honest answer to "what is knowledge?" includes a healthy list of things nobody has nailed down:

- **Can knowledge be analyzed at all?** Or was Williamson right that it's bedrock – a primitive we explain other things *with*, not *from*?

- **Internal or external?** Does being justified require reasons you can access by reflection, or just wiring that tends to produce truths?
- **One currency or two?** Is rational belief fundamentally graded (credence), all-or-nothing, or both somehow reconciled?
- **Is there really a universal human epistemology** – and if so, did *evolution* install the instinct that luck-based "knowing" doesn't count? (A thread for **Day 74**.)
- **Is the brain *literally* Bayesian**, or is "the brain does inference" just a useful way of describing it from outside?
- **And the question that will haunt the AI block:** when a system like the one that drafted this page outputs a true, well-supported claim, does it *know* anything – or is it the ultimate Gettier case, right for reasons that have nothing to do with the truth? (**Days 138–145**.)

## ◆ THE DAY IN THREE SENTENCES

## BIG IDEA

For 2,300 years knowledge looked like justified true belief — until Gettier showed in three pages that you can hold all three and still not know, because your reasons and the truth can meet by luck rather than by connection.

## BEST ANALOGY

The stopped clock that's right twice a day — and the archer whose arrow is blown off target, then blown back onto the bullseye: accurate, but not *apt*.

## LIVE CONTROVERSY

Whether the fix is a fourth condition (and which), whether knowledge is unanalyzable bedrock, and whether "belief" should give way to graded, Bayesian credence — with a real scientific frontier in the claim that the brain is a prediction machine.

---

THREADS TODAY > information (credence & the Bayesian brain) · energy (Friston's free energy) · computation (mind as inference engine) — with light first touches of emergence and evolution.

---

## — SOURCES

## Sources & further reading

1. Gettier, E. L. (1963). "Is Justified True Belief Knowledge?" *Analysis* 23(6): 121–123.  
doi:10.1093/analys/23.6.121. doi.org/10.1093/analys/23.6.121
2. Ichikawa, J. J. & Steup, M. "The Analysis of Knowledge." *Stanford Encyclopedia of Philosophy* (rev. 2018). plato.stanford.edu/entries/knowledge-analysis — JTB, the Gettier cases, safety/sensitivity, and the knowledge-first turn.

3. "Gettier problem." *Wikipedia* (accessed 2026). [en.wikipedia.org/wiki/Gettier\\_problem](https://en.wikipedia.org/wiki/Gettier_problem) – precedents in Russell (1948), Dharmottara (~770 CE), and Gaṅgeśa (14th c.).
4. Russell, B. (1948). *Human Knowledge: Its Scope and Limits*. London: Allen & Unwin. – the stopped-clock case (pp. ~170–171).
5. Goldman, A. (1976). "Discrimination and Perceptual Knowledge." *Journal of Philosophy* 73(20): 771–791. – the fake-barn case; reliabilism.
6. Nozick, R. (1981). *Philosophical Explanations*. Harvard University Press. – truth-tracking / sensitivity.
7. Sosa, E. (1999). "How to Defeat Opposition to Moore." *Philosophical Perspectives* 13: 141–153. – the safety condition. See also Sosa (2007), *A Virtue Epistemology* (apt belief).
8. Zagzebski, L. (1994). "The Inescapability of Gettier Problems." *The Philosophical Quarterly* 44(174): 65–73. – the recipe defeating any luck-excluding fix.
9. Williamson, T. (2000). *Knowledge and Its Limits*. Oxford University Press. overview – knowledge-first epistemology; knowledge as the most general factive mental state.
10. Weinberg, J. M., Nichols, S. & Stich, S. (2001). "Normativity and Epistemic Intuitions." *Philosophical Topics* 29(1–2): 429–460. – the founding (later contested) cross-cultural x-phi study.
11. Machery, E., Stich, S., Rose, D., Chatterjee, A., Karasawa, K., Struchiner, N., Sirker, S., Usui, N. & Hashimoto, T. (2017). "Gettier Across Cultures." *Noûs* 51(3): 645–664. doi:10.1111/nous.12110. doi.org/10.1111/nous.12110
12. Kim, M. & Yuan, Y. (2015). "No cross-cultural differences in the Gettier car case intuition: A replication study of Weinberg et al. 2001." *Episteme*. [philpapers.org/rec/KIMNCD](https://philpapers.org/rec/KIMNCD)
13. Weisberg, J. "Bayesian Epistemology." *Stanford Encyclopedia of Philosophy*. [plato.stanford.edu/entries/epistemology-bayesian](https://plato.stanford.edu/entries/epistemology-bayesian) – credences, conditionalization, and the Dutch book argument (Ramsey 1926; de Finetti 1937).
14. Clark, A. (2013). "Whatever next? Predictive brains, situated agents, and the future of cognitive science." *Behavioral and Brain Sciences* 36(3): 181–204. See also Clark, *Surfing Uncertainty* (OUP, 2016).
15. Friston, K. (2010). "The free-energy principle: a unified brain theory?" *Nature Reviews Neuroscience* 11(2): 127–138. doi:10.1038/nrn2787. doi.org/10.1038/nrn2787
16. Hohwy, J. (2013). *The Predictive Mind*. Oxford University Press.

## OPTIONAL APPENDIX

# Appendix: The Rest of the Map

*This section is optional supplemental reading. You can skip it without losing the main lesson.*

*We spent the main lesson on one belief, on one late morning. The field is far larger than one clock.*

**T**he main piece had a tight job: take a single belief – *it's 9:12* – and ask whether it counted as knowledge. To do that it leaned, quietly, on a stack of assumptions it never examined, and it strolled right past whole provinces of the subject without nodding. Does knowing require *certainty*? Can the skeptic who says you know *nothing* actually be answered? Does the word "know" even hold still from one sentence to the next? Why is knowledge worth *more* than a true belief that does the same job? And what about all the knowing that has nothing to do with facts – knowing how to swim, knowing a face, knowing a city? This appendix walks the rest of that map. Nothing here repeats the main lesson; it all hangs off its edges.

## ↪ CONTINUES DIRECTLY FROM

**Day 1 – What Is Knowledge?** There we built the three-legged stool (justified true belief), watched Gettier kick a leg out with three pages, toured the failed "fourth condition" patches, mapped Agrippa's trilemma, and ended at three frontiers: the cross-cultural test of "knowledge" intuitions, Bayesian credence, and the predictive brain. Keep two images from that day in your pocket – the *stopped clock* (right by luck, not connection) and the *archer* whose arrow is blown off course then back onto the bull (accurate, but not *apt*). Both come back transformed below.

## ◇ SEVEN ROOMS WE SKIPPED

1. **The trapdoors under Gettier** – the two hidden assumptions that make the trick possible, and the escape hatch (certainty) that drops you into skepticism.
2. **The skeptic at the door** – dreams, demons, brains in vats, and the 2020s simulation upgrade.

3. **"Knows" on a sliding scale** – the Bank Cases: same evidence, different stakes, opposite verdict.
4. **The luck we were really chasing** – anti-luck epistemology, which finally explains *why* the patch wars happened.
5. **Why knowing beats being right** – Meno's road, and the value of knowledge.
6. **The kinds of knowing we ignored** – how, and by acquaintance.
7. **Almost everything you know, someone told you** – testimony, disagreement, and epistemic injustice.

## §1 THE MACHINERY

### The two trapdoors under every Gettier case

Before we explore new rooms, look down. Gettier's three-page bomb only goes off because the floor has two trapdoors built into it – two assumptions so natural the main lesson never paused on them. Name them and the whole landscape reorganizes.

**Trapdoor one: fallible justification.** The classical picture lets you be *justified* in believing something that turns out *false*. Smith had excellent reason to believe "Jones will get the job" – the boss said so – and it was false. If justification had to *guarantee* truth, that step would be impossible and the case couldn't even start. **Trapdoor two: closure.** Justification (and knowledge) is assumed to travel across *entailment*: if you're justified in believing something, you're justified in believing what it obviously implies. Smith reasons from "Jones will get it (and has ten coins)" to the weaker "the winner has ten coins" – a valid inference – and carries his justification along for the ride. Knock out either plank and Gettier cases evaporate.

That hands us a tempting exit. Slam trapdoor one shut: insist that real knowledge needs *infallible* justification – reasons that make error literally impossible. No more Gettier cases, ever. This is the dream of *infallibilism*, and it is very old. Descartes went looking in 1641 for a single belief no demon could fake, and found exactly one that survives even the supposition that an all-powerful deceiver is fooling you about everything else: *I think, therefore I am*. You cannot be tricked into wrongly believing you exist, because the tricking requires a you to be tricked.

The trouble is what the demon takes with him on the way out. If knowledge demands that kind of certainty, then you do not know you have hands, that the sun will rise, that the

person across the table is your friend and not an android – because a clever enough deception could fake any of it. Buy certainty and the price is **skepticism**: the bar is set so high that almost nothing clears it. Peter Unger argued exactly this in *Ignorance* (1975) – that "knows," used strictly, applies to virtually nothing, much as "flat" strictly applies to no real surface. So infallibilism doesn't dissolve the problem; it trades a small puzzle (the odd lucky belief) for a total one (you know next to nothing). Which is our cue to open the next door, where that skeptic is already knocking.

#### GETTIER'S OTHER CASE, IN ONE BREATH

The main lesson used the coins. Gettier's *second* case shows trapdoor two even more nakedly. Smith, with great evidence, believes "Jones owns a Ford." From that he validly deduces "Jones owns a Ford, *or* Brown is in Barcelona" – a disjunction he's entitled to, since a true disjunct makes the whole thing true. But Jones doesn't own a Ford after all... and Brown, by pure fluke, *is* in Barcelona. The disjunction is true, justified, believed – and obviously not known. Closure carried the justification; luck supplied the truth. Same skeleton, fancier clothes.

### — §2 THE BIGGEST OMISSION

## The skeptic at the door

Western epistemology has a recurring houseguest who refuses to leave: the figure who says you can't know *anything* about the world outside your own mind. The main lesson kept the door shut. Open it, because every modern theory of knowledge is partly built to deal with what's standing there.

The skeptic's tools are thought experiments, escalating in cruelty. First, the **dream**: right now, how do you know you're not asleep? Dreams feel utterly real from the inside; you've been fooled before. (The Daoist Zhuangzi, around 300 BCE, dreamt he was a butterfly and woke unsure whether he was a man who'd dreamt a butterfly or a butterfly now dreaming a man – the same wound the Buddhist Dharmottara reopened in the main lesson, proof again that minds keep tripping over this independently.) Descartes raised the stakes to an **evil demon** bent on deceiving you about everything. The twentieth century updated the hardware: you might be a *brain in a vat*, nerves wired to a computer feeding you exactly the experiences you're having now (Hilary Putnam, *Reason, Truth and History*, 1981). You cannot tell from the inside. That's the whole point.

Spelled out, the skeptic's argument is brutally clean – and it runs on the very closure principle from §1:

*(1) You don't know you're not a handless brain in a vat being fed a hand-experience.*

*(2) If you know you have hands, then (since having hands entails not being a handless vat-brain) you know you're not one.*

*(3) So you don't know you have hands.*

Each line looks reasonable; together they seem to prove you know nothing about the external world. The interactive below lets you try every way out – and discover that each "way out" is a named philosophical position with a price tag.

### The Skeptic's Syllogism, as four exits

MOVE	LINE REFUSED	REPRESENTATIVE VIEW	COST
Accept all three	None	Skepticism	You do not know you have hands, or much about the external world.
Reject P1	You do not know you are not a vat-brain	Moore's common-sense reply	Can feel like insisting rather than explaining.
Reject P2	Closure	Dretske / Nozick relevant alternatives	Closure is deeply intuitive and useful elsewhere.
Change the standard	A fixed meaning of "know"	Contextualism	The skeptic wins in the seminar; ordinary speakers win in ordinary life.

The doors are worth naming in full. **G. E. Moore** (1939) simply ran the argument backwards: *I am far more sure that here is one hand (holding it up) than I am of any fancy premise the skeptic offers* – so if the premises imply I don't know it, so much the worse for the premises. Cheeky, and strangely hard to beat. **Fred Dretske** (1970) and Robert Nozick (1981) took the surgical route: *deny closure*. On Dretske's *relevant alternatives* view, to know something you only need to rule out the *relevant* ways you could be wrong, not every bizarre one. At the zoo you know the animal is a zebra – you've ruled out "it's a horse," "it's a goat" – even though you haven't ruled out "it's a mule cleverly painted to look like a zebra," because in this context that's not a live possibility. Knowledge doesn't automatically transmit to every entailment. The cost is steep: closure is intuitive, and giving it up has consequences elsewhere. **Contextualism** (our next section) offers the diplomat's solution: maybe the skeptic and Moore are *both* right, because "know" means something stricter in the skeptic's seminar than in ordinary life.

## The 2020s upgrade: are we in a simulation?

The vat got a software update. Nick Bostrom's **simulation argument** (*Philosophical Quarterly*, 2003) makes a careful probabilistic case that at least one of three things is true: civilizations almost never reach the technology to run ancestor-simulations; or they reach it but choose not to; or *we are almost certainly living in one*. David Chalmers, in **Reality+** (2022), takes the next step and bites a bullet most people won't: he argues we *can't know* we're not simulated and should assign the possibility real probability – but that this **isn't a catastrophe**, because *"virtual reality is genuine reality."* A simulated tree, on his *simulation realism*, is a real digital object, not an illusion; if you've always lived in a perfect simulation, your belief "that's a tree" is *true*, just realized in silicon. The skeptic assumed a fake world means false beliefs; Chalmers denies the link.

Two honest labels before we move on. The simulation *hypothesis* – that we are in fact simulated – is, as it stands, **untestable metaphysics, not science**: there's no agreed observation that would confirm or refute it, which puts it on the wrong side of the demarcation line we'll draw tomorrow. [UNFALSIFIABLE] The *philosophical* payoff is real all the same: it sharpens what we even mean by "real" and "know." And there's a famous reply that turns the screw the other way. Putnam argued that "I am a brain in a vat" is **self-refuting**: your words only mean what they do because of your causal history, so a lifelong brain-in-a-vat's word "vat" couldn't refer to real vats (it never causally touched one) – meaning that if you *were* a vat-brain, your sentence "I am a brain in a vat" would come out *false*. Whether that works is still argued, which is precisely why this thread runs straight into the AI block: when a system trained only on text outputs "Paris is in France," does it *know* that – or is it

the purest brain-in-a-vat of all, with words that never touched the world? Hold the question for **Days 138–145**.

— §3 THE MOVING TARGET

## "Knows" on a sliding scale

Here is a possibility the main lesson never entertained: maybe sixty years of hunting the perfect definition of "knows" failed because the word was never aiming at a fixed point. Consider a pair of cases from Keith DeRose (*Philosophy and Phenomenological Research*, 1992) that have launched a thousand papers – the **Bank Cases**.

It's Friday. You drive past your bank, which has a long Saturday line, and decide to come back tomorrow. Your wife asks if it'll be open Saturday. *Low stakes* version: nothing much rides on it; you say, "Yes, I know it's open Saturdays – I was here two Saturdays ago." That sounds true. You know it. *High stakes* version: there's a check that *must* be deposited by Monday or you bounce your mortgage and lose the house, and your wife points out, reasonably, that banks do change their hours. Now the very same sentence – "I know it's open Saturday" – curdles in your mouth. "Well... I'd better go in and check." Same person, same memory, same evidence, same day. Only the stakes (and whether someone raised the chance of error) have changed. Yet the knowledge seems to come and go. The dial below lets you slide between the two and watch it flip.

## The Bank Cases, as a stakes table

CASE	EVIDENCE	STAKES	NATURAL VERDICT	WHAT IT TESTS
Low stakes	You were there two Saturdays ago.	A minor errand.	"I know it is open."	Ordinary standards are easy to meet.
High stakes	The same memory.	A mortgage deadline.	"I had better check."	Whether practical stakes affect knowledge.
Error raised	The same memory plus a live doubt.	Any serious consequence.	The claim to know weakens.	Whether context shifts the word or the knower's state.

Three camps, three diagnoses of the same data. **Contextualism** (DeRose; David Lewis, "Elusive Knowledge," 1996; Stewart Cohen, 1988) locates the shift in the *word*: "knows" is like "tall" or "here" – context-sensitive. Raising the stakes or mentioning error raises the standard a belief must meet for the sentence "S knows" to count as true. Both utterances are correct, in their own conversations. The skeptic is even right in the seminar – he's just jacked the standard sky-high. **Pragmatic encroachment** (Jason Stanley, *Knowledge and Practical Interests*, 2005; Fantl & McGrath; John Hawthorne, *Knowledge and Lotteries*, 2004) puts the shift in the *knower*: what *you* know genuinely depends on what's practically at stake *for you*, because knowledge is supposed to be the thing you can act on. High stakes really can deprive you of knowledge you'd have had when it didn't matter – a startling idea, since it lets practical pressure "encroach" on a supposedly purely factual state. **Invariantism** (the traditional holdout) digs in: "knows" means one fixed thing, the standards don't move, and one of your two verdicts is simply mistaken – you either knew all along or never did, and the stakes just changed how *willing* you were to *say* so. [AGREED] [UNRESOLVED] The data is robust; its explanation is one of the most active fault lines in the field.

## §4 THE PATTERN BEHIND THE PATCHES

## The luck we were really chasing

Return to the patch wars from the main lesson – no-false-lemmas, sensitivity, safety, virtue. They looked like a grab-bag of clever fixes that each met a nastier counterexample. Step back and they snap into focus: every one was chasing the *same ghost*. Duncan Pritchard gave it a precise name in *Epistemic Luck* (Oxford, 2005). The enemy of knowledge is a specific species he calls *veritic luck*: your belief is true in the actual world, but in *almost all the nearby ways things could have gone*, you'd have believed the same thing and been wrong. The truth and your believing it are only accidentally in step.

This is the deep content of the "safety" idea, and it's worth *seeing*. Picture the actual world as a dot, ringed by the nearby possible worlds – the small, realistic variations on how things might have been. A belief is *safe* (knowledge-grade) when it stays true across that neighborhood, and *unsafe* (merely lucky) when a slight nudge flips it to false. Toggle the three scenarios below and watch the neighborhood light up.

### Safe vs. Lucky, as nearby-worlds cases

SCENARIO	ACTUAL WORLD	NEARBY WORLDS	VERDICT
Working clock	Your belief is true.	Small variations still leave you right.	Safe: knowledge-grade.
Stopped clock	Your belief is true at 9:12.	A minute earlier or later, the same belief is false.	Unsafe: veritic luck.
Fake-barn country	You see the one real barn.	Most nearby looks would have landed on facades.	Unsafe: environmental luck.

That single picture retroactively explains the whole mess. The stopped clock fails *hard* – a minute either side and you're wrong, so the neighborhood is a sea of red. Fake-barn country is subtler: the barn you're looking at is genuinely there (the core is green), but you're

surrounded by façades, so a glance a hundred meters either way would have fooled you – red neighborhood, no knowledge, even with a true justified belief and no false premise. The patches all failed because each tried to capture "green neighborhood" with a slightly different yardstick, and luck kept finding the gaps.

Two more patches the main lesson didn't name, now that we have the frame. **Defeasibility theory** (Lehrer & Paxson, 1969) said knowledge is *un-defeated* justified true belief: there must be no true fact out there that, if you learned it, would dissolve your justification. It handles many cases elegantly – until the "misleading defeater" twist, where there's a true-but-misleading fact that *shouldn't* rob you of knowledge but technically does, forcing ever-finer distinctions. And reaching back further, the **causal theory** (Goldman, 1967, before he turned reliabilist) demanded that the fact *cause* your belief – no causal chain, no knowledge. Beautiful for perception; fatal for mathematics, since the number 7 and the Pythagorean theorem don't cause anything (Paul Benacerraf pressed exactly this "access problem" in 1973). You can't shake hands with an abstract object.

And the deepest crack in reliabilism, which the main lesson only gestured at: the **generality problem** (Conee & Feldman, 1998). Reliabilism says a belief is justified if produced by a *reliable process* – but *which* process? Your belief that it's 9:12 was produced by "reading a clock," and also by "reading *that* clock," and "using vision in dim light," and "trusting instruments on Tuesdays" – each as real as the others, each with a different reliability score. Pick the type and you've picked the verdict. Specifying the "right" grain, in a principled way, has proven stubbornly hard.

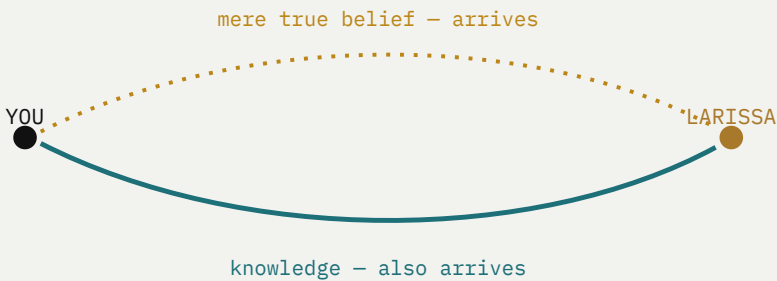
Where does Pritchard land? At *anti-luck virtue epistemology*: knowledge needs *both* conditions, because they catch different failures. You need **safety** (a green neighborhood – no veritic luck) *and* you need **aptness** (the belief is true *through your own ability* – the archer's skill from the main lesson). Neither alone suffices: the stopped clock can fail safety, fake-barn country can have local skill but bad luck. It's not a tidy three-word formula – and that, by now, may be the lesson. Knowledge might just *be* the kind of thing that takes two independent guarantees, one about you and one about your world.

## — §5 THE QUESTION UNDER THE QUESTION

### Why is knowing worth more than just being right?

Step back from "what is knowledge?" to a question Plato asked first and nobody has fully answered: *why do we care?* If a true belief gets the job done, what does the extra machinery of knowledge buy you? Plato put it as a traveler's problem in the *Meno* (~380 BCE).

Suppose you want to walk to the town of Larissa. A person who *knows* the road will get you there. But so will a person who merely has a *true belief* about the road – who's never been, but happens to be right. For the purpose of arriving, the two are worth exactly the same. So why has the entire tradition prized knowledge above true belief? This is the *value problem*, and it's a load-bearing question: a theory of knowledge that can't say why knowledge is *better* has arguably missed the point of the concept.



If both roads reach Larissa, what is the second one worth?

The value problem turns into a precise weapon against one of the main lesson's theories. It's called the *swamping problem* (Linda Zagzebski, 2003). Reliabilism says knowledge is true belief from a reliable process. But ask *what the reliability adds in value*. Reliability is good only because it tends to produce truth. So once you *already have* the truth, what does it add that this particular truth also came from a reliable source? Zagzebski's homely analogy: a cup of good coffee is no *better* to drink for having come from a reliable coffee machine rather than an unreliable one that happened to produce an identical cup. The good-making feature (deliciousness / truth) is already present; the source's reliability gets *swamped*, adding nothing. If that's right, reliabilism can't explain why knowledge beats lucky true belief – the very thing a theory of knowledge most needs to deliver.

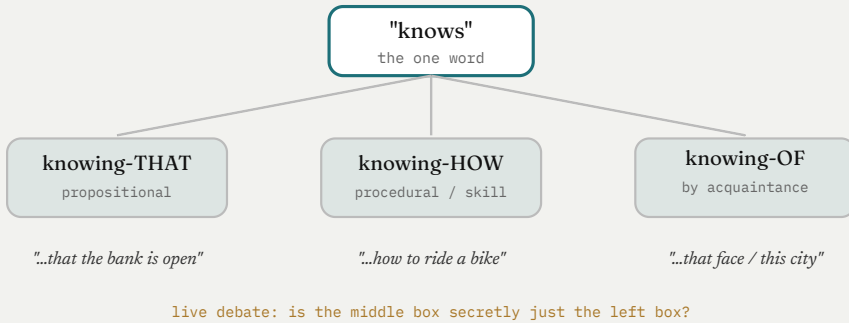
This is where **virtue epistemology** earns its keep, and where the archer finally pays off. Its answer: knowledge isn't valuable as a *better-stocked* true belief; it's valuable as an *achievement* – a success that's *yours*, that came about *through your own competence*. And achievements carry a kind of worth that lucky successes never do, the way a bullseye you actually aimed is worth something a lucky gust-blown hit isn't, even though the arrow lands in the same spot. A true belief reached through your own cognitive skill is a *cognitive achievement*; a lucky true belief is not. That's the extra value – not in the result, but in the

*getting there.* The road to Larissa you can actually find again is worth more than the one you stumbled onto, even on a day you both arrive.

## — §6 THE OTHER KNOWINGS

### The kinds of knowing we ignored

Everything so far – the entire main lesson – was about *propositional knowledge*, knowledge-*that*: knowing *that* it's 9:12, *that* Jones got the job. But look how much of "know" in plain English isn't that at all. You know *how* to ride a bicycle. You know your mother's face. You know Lisbon. None of these is a stockpile of facts, and philosophers have argued for a century about how they relate.



One English verb, at least three different relations to the world.

**Knowing-how.** Gilbert Ryle, in *The Concept of Mind* (1949), insisted that knowing how to do something is not knowing a set of facts. A brilliant cyclist may be unable to state a single law of balance; a person who has memorized every fact about bicycles may topple on the first try. Worse, Ryle argued, reducing skill to facts triggers a regress: if every skilled act required first *knowing the proposition* describing the rule, you'd then need the skill of *applying* that rule, which would need another rule, forever. So skill must be its own kind of knowing. The twist: Jason Stanley and Timothy Williamson fired back in "**Knowing How**" (2001) with *intellectualism* – the claim that knowing-how just *is* a species of knowing-that after all (knowing, of some way to ride, *that* it is a way to ride), dressed in different grammar. Whether skill collapses into propositions is genuinely unsettled. [CONTESTED]

**Knowing by acquaintance.** Bertrand Russell (1911) drew a second cut: between knowledge *by acquaintance* – your direct, unmediated grip on a patch of red you're seeing, a pain you're feeling, a face you're looking at – and knowledge *by description*, the facts you know *about* things you've never directly met ("the first person to stand on the Moon," whom you know only as the one satisfying that description). You can know a stupendous amount *about* Bismarck and never have known *him*; you know the color red in a way the world's greatest blind physicist, who knows every fact about wavelengths, does not. That gap – facts about an experience versus the experience itself – is a quiet seed for the hardest problem in the entire course, the one waiting on **Day 123**: why there's *something it is like* to see red at all.

## — §7 THE SOCIAL TURN

### Almost everything you know, someone told you

The main lesson, like most of traditional epistemology, imagined a lone mind facing the world – one person, one clock. But run an audit of what you actually know. That the Earth is about 4.5 billion years old. That Antarctica exists. Your own date of birth. The boiling point of water. You verified essentially none of it first-hand; you were *told*, by teachers, books, parents, instruments, strangers. *Testimony* is the overwhelming bulk of any human being's knowledge – and for centuries epistemology treated it as an afterthought.

The central question is whether trusting testimony is something you have to *earn* or something you're *entitled* to by default. **David Hume** (1748) took the demanding line: testimony is only as good as your own inductive track record of when testimony has proved reliable – it *reduces* to evidence you've personally gathered. **Thomas Reid** (1764) found this absurd: no child could bootstrap a track record before trusting anyone, and in fact we're built with a "principle of credulity," a default disposition to believe what we're told, exactly as we're built to trust our senses. On Reid's *anti-reductionist* view, testimony is a *basic* source of knowledge, not a derived one – and it has to be, or knowledge couldn't get off the ground in a social animal. The modern field mostly agrees that some default trust is unavoidable; the fights are over how much, and when it's defeated.

Two newer rooms branch off this one, and both matter enormously in 2026. The first is **disagreement**. When someone you regard as an *epistemic peer* – as smart, as informed, as careful as you – looks at the same evidence and concludes the opposite, what should you do? The *conciliationist* or "equal-weight" view (Adam Elga, *Noûs*, 2007; David Christensen, 2007) argues you should move substantially toward them: to stay put is to claim, with no independent reason, that *you're* the one who got it right and they made the mistake. The *steadfast* view answers that sometimes you can rationally hold your ground, because your

own reasoning is evidence too. It sounds abstract until you notice it's the whole epistemology of echo chambers, expert consensus, and what to do when half your sources contradict the other half. [DEBATE]

The second is sharper still: **epistemic injustice**, named by Miranda Fricker (*Epistemic Injustice: Power and the Ethics of Knowing*, 2007). Because so much knowing runs on testimony, *who gets believed* becomes an ethical question, not just an epistemic one. Fricker isolates two wrongs. *Testimonial injustice*: a speaker's word is given less credence than it deserves because of prejudice about who they are – the patient whose pain is dismissed, the witness disbelieved for their accent or gender. *Hermeneutical injustice*: subtler and deeper – a person can't even make sense of their own experience, to themselves or others, because the surrounding culture hasn't yet developed the *concept* for it (her example: the experience we now call sexual harassment, suffered by people who had no word for it and so couldn't name the wrong). Knowledge, it turns out, has a politics: the tools for understanding are unevenly distributed, and that unevenness can itself be an injustice.

#### THE FUNCTION-FIRST ESCAPE HATCH

There's a radical way to end the whole 180-page hunt for a definition, and it threads the social turn back to the start. Edward Craig, in *Knowledge and the State of Nature* (1990), proposed: stop asking "*what is knowledge?*" and ask "*what is the concept FOR – why would creatures like us ever invent it?*" His answer: a social, language-using species desperately needs a way to flag **good informants** – to mark out whose word you can act on. "Knowledge" is the tag we evolved to pin on reliable sources of true information. That instantly explains the things the analyses struggled with: why knowledge must be *true* (a tip that's false is worthless), why *luck* disqualifies (you can't rely on a fluke next time), and why we care at all (survival in a world where most of what you need to know, you must get from others). It rhymes with Williamson's "stop trying to define it," and it cashes out the main lesson's open question – did *evolution* install the instinct that luck-based knowing doesn't count? Craig's answer is essentially: yes, and here's why it would.

#### — § 8 THE FORMAL EDGE

## Two more frontiers, beyond Bayes

The main lesson's formal frontier was Bayesian credence. Two further formal ideas deserve a place on the map, because both bite ordinary intuitions and both feed straight into computer science and AI.

**The logic of knowing.** Jaakko Hintikka, in *Knowledge and Belief* (1962), treated "knows" as a formal operator you can reason with, like "necessarily" – launching *epistemic logic*, now a workhorse in computer science (reasoning about what distributed agents and AI systems "know"). It immediately surfaces deep puzzles. The *KK principle*: if you know  $p$ , do you thereby know *that you know  $p$* ? Tempting, but Williamson (from the main lesson) argues it's false – you can know something without being in a position to know that you know it, because knowledge has blurry margins. And *logical omniscience*: the clean logic implies that if you know some axioms, you know *every* logical consequence of them – which would make every mathematician instantly aware of every theorem. Obviously false for real, bounded minds, and a central headache for modeling actual reasoners (and machines).

**The preface paradox.** A companion to the lottery paradox from the main lesson, and arguably nastier. You write a long, careful book. For *each* claim in it, you've checked your work and rationally believe it's true. Yet you also write, sincerely, in the preface: "no doubt errors remain, and they are mine alone" – because you know that across hundreds of claims, you've almost certainly slipped *somewhere*. So you rationally believe each individual claim, and *also* rationally believe that *at least one of them is false* (David Makinson, "The Paradox of the Preface," 1965). Those can't all be true together. The moral lands on the main lesson's open question with full force: ordinary all-or-nothing belief isn't *closed under conjunction* – believing each of many things doesn't license believing their grand conjunction – which is one more reason the field keeps drifting from yes/no belief toward graded credence. The dial, again, doing what the switch can't.

## ◆ THE APPENDIX IN THREE SENTENCES

## BIG IDEA

The main lesson made knowledge look like one tidy puzzle — find the fourth condition — but it's really a constellation: whether certainty is required (and the skeptic that demand invites), whether "knows" even holds still as stakes change, what knowledge is *worth* over mere true belief, and the fact that almost all of it comes from *other people*.

## BEST NEW ANALOGY

The neighborhood of nearby possible worlds: knowledge is a belief whose neighborhood stays green (safe), while luck is a belief one nudge from red — and the road to Larissa you can find *again* is worth more than the one you stumbled onto, even when both arrive.

## LIVE CONTROVERSY

Why the Bank-Case verdict flips — context shifting the *word* "knows" (contextualism), stakes shifting what the *knower* knows (pragmatic encroachment), or neither (invariantism) — is among the field's hottest open fault lines, alongside whether closure can be denied and whether knowing-how is secretly knowing-that.

---

THREADS HERE > information (testimony & the social transmission of knowledge; preface/credence) · computation (epistemic logic; modal "neighborhoods" of worlds) · evolution (Craig: the concept of knowledge as a good-informant detector built for a social species) — picking up the same five we're tracking all 180 days.

## — OPEN QUESTIONS

## What this appendix leaves unsettled

- **Certainty or not?** Is the infallibilist right that real knowledge needs error-proof reasons (inviting skepticism) – or is fallible knowledge the only kind worth wanting?
  - **Can closure be denied without disaster?** Dretske and Nozick block the skeptic by giving it up; the cost elsewhere is still being counted.
  - **Does "knows" move?** Context-sensitive, stakes-sensitive, or fixed – and if it moves, what exactly is moving, the word or the world?
  - **Can the value of knowledge be explained at all,** or does every account leave knowledge looking no better than lucky true belief?
  - **Is knowing-how just knowing-that** in disguise, or its own irreducible kind of grip on the world?
  - **Is testimony basic or earned?** – and, downstream, when a peer disagrees, must you really meet them halfway?
  - **And the function-first wager:** if the concept of knowledge exists to flag good informants, does that *dissolve* the analysis project – or just relocate it?
- 

## — SOURCES

## Sources & further reading

Classical works are cited by original date; all are standard, widely available editions. Verified secondary anchors and reference entries are linked.

1. Descartes, R. (1641). *Meditations on First Philosophy*. – methodic doubt, the evil demon, and the cogito as the one indubitable point.
2. Unger, P. (1975). *Ignorance: A Case for Scepticism*. Oxford University Press. – infallibilism pushed to its skeptical conclusion ("knows," like "flat," applies to almost nothing).
3. Moore, G. E. (1939). "Proof of an External World." *Proceedings of the British Academy* 25: 273–300. – "Here is one hand": running the skeptical argument in reverse.
4. Dretske, F. (1970). "Epistemic Operators." *Journal of Philosophy* 67(24): 1007–1023. – denying closure; the relevant-alternatives view; the zebra/painted-mule case.

5. Nozick, R. (1981). *Philosophical Explanations*. Harvard University Press. – sensitivity / truth-tracking and its own denial of closure.
6. Putnam, H. (1981). *Reason, Truth and History*. Cambridge University Press. – the brain-in-a-vat, and the semantic-externalist argument that "I am a BIV" is self-refuting.
7. Bostrom, N. (2003). "Are You Living in a Computer Simulation?" *Philosophical Quarterly* 53(211): 243–255. [simulation-argument.com](http://simulation-argument.com)
8. Chalmers, D. J. (2022). *Reality+: Virtual Worlds and the Problems of Philosophy*. W. W. Norton / Allen Lane. – "virtual reality is genuine reality"; simulation realism. [consc.net/reality](http://consc.net/reality)
9. DeRose, K. (1992). "Contextualism and Knowledge Attributions." *Philosophy and Phenomenological Research* 52(4): 913–929. – the Bank Cases. See also DeRose (1995), "Solving the Skeptical Puzzle," *Philosophical Review* 104(1): 1–52.
10. Lewis, D. (1996). "Elusive Knowledge." *Australasian Journal of Philosophy* 74(4): 549–567. – contextualism and the rule of attention.
11. Cohen, S. (1988). "How to Be a Fallibilist." *Philosophical Perspectives* 2: 91–123. – the airport cases.
12. Stanley, J. (2005). *Knowledge and Practical Interests*. Oxford University Press. – pragmatic encroachment / interest-relative invariantism. See also Hawthorne, J. (2004), *Knowledge and Lotteries* (OUP); Fantl, J. & McGrath, M. (2009), *Knowledge in an Uncertain World* (OUP).
13. Pritchard, D. (2005). *Epistemic Luck*. Oxford University Press. – the modal account of luck; veritic luck; the safety condition; later, anti-luck virtue epistemology. Overview: IEP, "Epistemic Luck."
14. Lehrer, K. & Paxson, T. (1969). "Knowledge: Undefeated Justified True Belief." *Journal of Philosophy* 66(8): 225–237. – the defeasibility analysis.
15. Goldman, A. (1967). "A Causal Theory of Knowing." *Journal of Philosophy* 64(12): 357–372. – and Benacerraf, P. (1973), "Mathematical Truth," *J. Phil.* 70(19): 661–679, on why it fails for abstract objects.
16. Conee, E. & Feldman, R. (1998). "The Generality Problem for Reliabilism." *Philosophical Studies* 89(1): 1–29.
17. Plato. *Meno* (~380 BCE). – the road to Larissa; the value problem (knowledge vs. true belief).
18. Zagzebski, L. (2003). "The Search for the Source of Epistemic Good." *Metaphilosophy* 34(1–2): 12–28. – the swamping problem. See also Kvanvig, J. (2003), *The Value of Knowledge and the Pursuit of Understanding* (Cambridge UP).
19. Ryle, G. (1949). *The Concept of Mind*. University of Chicago Press. – knowing-how vs. knowing-that; the regress of rules.
20. Stanley, J. & Williamson, T. (2001). "Knowing How." *Journal of Philosophy* 98(8): 411–444. – intellectualism: knowing-how as a species of knowing-that.

21. Russell, B. (1910–11). "Knowledge by Acquaintance and Knowledge by Description." *Proceedings of the Aristotelian Society* 11: 108–128.
22. Hume, D. (1748). *An Enquiry Concerning Human Understanding*, §X. – the reductionist view of testimony. Reid, T. (1764). *An Inquiry into the Human Mind on the Principles of Common Sense*. – testimony as a basic source (anti-reductionism).
23. Elga, A. (2007). "Reflection and Disagreement." *Noûs* 41(3): 478–502. doi:10.1111/j.1468-0068.2007.00656.x. And Christensen, D. (2007), "Epistemology of Disagreement: The Good News," *Philosophical Review* 116(2): 187–217.
24. Fricker, M. (2007). *Epistemic Injustice: Power and the Ethics of Knowing*. Oxford University Press. – testimonial and hermeneutical injustice.
25. Craig, E. (1990). *Knowledge and the State of Nature: An Essay in Conceptual Synthesis*. Oxford University Press. – the function-first / good-informant genealogy of the concept.
26. Hintikka, J. (1962). *Knowledge and Belief: An Introduction to the Logic of the Two Notions*. Cornell University Press. – epistemic logic; the KK principle; logical omniscience.
27. Makinson, D. C. (1965). "The Paradox of the Preface." *Analysis* 25(6): 205–207.
28. Reference surveys: *Stanford Encyclopedia of Philosophy* – "Skepticism," "Epistemic Contextualism," "The Value of Knowledge," "Epistemological Problems of Testimony," "Epistemic Injustice."

## OPTIONAL APPENDIX

# Appendix: The Edge of the Map

*This section is optional supplemental reading. You can skip it without losing the main lesson.*

*The coastline still being drawn. Recent, high-voltage, and – every line of it – not yet safe to stand on.*

**T**he first appendix charted the *settled* hinterland – skepticism, contextual "knows," the value of knowing, the social web – territory mapped decades or centuries ago. This one sails to the edge, where the cartographers are still arguing about where the shore is. Everything below is peer-reviewed work from **2020 onward** that could genuinely redraw what we mean by "knowledge" – and precisely *because* it's that new, the hype filter does the heavy lifting. Nothing here is bankable. Each frontier comes with its own counter-literature already forming, and every claim wears a tag: [ESTABLISHED] [PROMISING] [CONTESTED] Read it the way you'd read a dispatch from an expedition still underway – thrilling, partial, and subject to revision by the next ship back.

## ↪ THIRD IN A SEQUENCE

**Day 1 – What Is Knowledge?** built the stool and watched Gettier kick a leg out. **Appendix I – "The Rest of the Map"** walked the settled provinces: skepticism, contextualism, anti-luck epistemology, the value problem, testimony and epistemic injustice. This piece is the live edge of that same continent. Where the social-turn rooms of Appendix I (testimony, disagreement, who gets believed) described the *structure* of knowing-from-others, several frontiers here describe what happens when that structure comes under deliberate *attack* – by manipulators, by machines, by the feed.

## ◆ SIX PLACES THE SHORELINE IS MOVING

1. **The zetetic turn** – epistemology pivots from *belief-states* to the *act of inquiry*, and finds its old rules in conflict with the new ones.

2. **Knowledge before belief** – cognitive science flips the furniture: maybe representing *knowledge* is more basic than representing belief.
3. **Do machines know – or bullshit?** – the philosophy of large language models, and a deliberately rude diagnosis.
4. **The epistemic backstop collapses** – deepfakes quietly remove a support that's been holding up testimony all along.
5. **Hostile epistemology** – echo chambers, manufactured clarity, and trust as something that can be *weaponized*.
6. **Accuracy-first** – a formal refoundation that re-derives Day 1's Dutch book from *truth* instead of *money*.

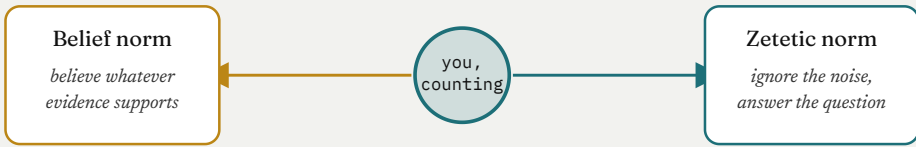
## §1 THE ACTIVITY TURN

# Epistemology forgot to ask about *inquiry*

[PROMISING] [CONTESTED]

Here is a strange omission, and once you see it you can't unsee it. For a century, epistemology has been almost entirely a theory of *states* – of belief, justification, knowledge: snapshots of a mind that has already finished thinking. It has had remarkably little to say about the *activity* that produces those states – the messy business of *inquiry*: asking a question, deciding what evidence to gather, knowing when to stop. Jane Friedman put a name to the missing half and lit a fuse under the field. The norms of inquiry she calls *zetetic* norms (from the Greek *zêtein*, to seek), and her landmark paper "The Epistemic and the Zetetic" (*Philosophical Review*, 2020) argues something genuinely destabilizing: the norms of *inquiry* and the norms of *belief* are not merely separate – they actively **conflict**.

The engine is a principle so obvious it sounds like a truism – the *Zetetic Instrumental Principle*: if you want to figure out the answer to a question, you ought to take the necessary means to figure it out. Now watch it collide with a bedrock epistemic norm – the evidentialist's command that you may believe whatever your evidence already supports. Suppose you're trying to count the windows on the building across the street. Good inquiry says: *focus, go count the windows, don't get distracted*. But at every passing instant your senses hand you sufficient evidence to form, and be permitted to believe, a thousand idle truths – the color of that car, the number of people on the corner, the shape of a cloud. The belief-norm *permits* all of them. The inquiry-norm tells you to *ignore* all of them and count windows. Conform to one and you flout the other. The diagram makes the squeeze concrete.



conform to one → violate the other

Friedman's tension: good inquiry and permissible belief pull opposite ways.

Why does this matter beyond the seminar? Because it suggests epistemology has been studying the wrong unit. If belief-norms and inquiry-norms genuinely clash, then a theory built only on belief is incomplete – maybe even *backwards*. The radical proposal, the "zetetic turn," is that **all epistemic norms are ultimately norms of inquiry** (suspension of judgment becomes a *question-directed* attitude; believing an answer is a way of *closing* a question). The field has not swallowed this whole – and that's the honest part. Arianna Falbo ("Should epistemology take the zetetic turn?", *Philosophical Studies*, 2023) and others argue the inquiry-norms are really *practical*, not distinctively epistemic, and that a purely zetetic epistemology can't explain why some beliefs are irrational even when believing them would *help* your inquiry. So: the *puzzle* is now taken with great seriousness across the field; the *grand thesis* that inquiry swallows everything is a genuine, unresolved fight. Either way, the question "what is knowledge?" is quietly being reframed as "what is it to inquire *well*?" – and that reframing reaches all the way forward to **Day 2**, where the scientific method is exactly a set of norms for collective inquiry.

## §2 THE COGNITIVE-SCIENCE TURN

### What if knowledge comes *before* belief?

[PROMISING] [CONTESTED]

Day 1 treated knowledge as something *built up from* belief: take a belief, add truth, add justification, screen out luck. Almost every theory we met assumed belief is the raw material and knowledge the finished product. A large interdisciplinary team led by Jonathan Phillips and Joshua Knobe dropped a target article in *Behavioral and Brain Sciences* – "Knowledge

before belief" (2021) – arguing that, as a matter of how human (and animal) minds actually work, this may be exactly **upside down**.

The standard story in psychology is that our "theory of mind" – our capacity to model other minds – is centered on *belief*, and matures when a child finally passes the *false-belief test* (understanding that someone can hold a belief the child knows to be false) around age four. Phillips and colleagues marshal converging evidence that representing *knowledge* is the more basic feat. The threads: developmentally, infants and toddlers track who has *seen* or has *access to* information – who *knows* – well before they can handle false beliefs. Comparatively, non-human great apes show robust signs of tracking what others can perceive and know, while convincing evidence of belief-tracking remains elusive. And in adults, attributions of knowledge are made at least as fast as – sometimes faster than – attributions of belief, which is hard to explain if "X knows p" is computed by first building "X believes p" and then checking extra conditions. The picture they paint is a *factive theory of mind*: the mind's first and most basic tool for modeling others is "what do they know?"; with the trickier, error-tolerant "what do they merely *believe* (perhaps falsely)?" coming later and costing more.



The proposed ordering – the reverse of the textbook "belief-first" story.

If it holds, the payoff for Day 1 is direct and large. It would be empirical ammunition for Timothy Williamson's *knowledge-first* program – the philosophical claim (which we met as pure armchair theory) that knowledge is the basic, unanalyzable state and belief should be explained in terms of *it*, not the reverse. Suddenly that's not just a logician's hunch; it's a candidate fact about the architecture of cognition, with an evolutionary rationale (a social animal urgently needs to track *who has reliable information* – echoing Edward Craig's "good-informant" genealogy from Appendix I). But honesty demands the asterisk, and here it's loud: a BBS target article arrives wrapped in dozens of peer commentaries, and many dissent hard. Critics argue the knowledge/belief line is blurrier than the authors allow, that "tracking who saw what" needn't be a full representation of *knowledge*, and that culture and

language shape the whole picture. So: the *behavioral findings* (kids and apes track informational access early) are reasonably solid; the *strong interpretation* (knowledge-representation is metaphysically and computationally prior, with belief assembled from it) is very much live. The furniture of the mind is being rearranged in real time, and the movers don't yet agree on the floor plan.

— §3 THE MACHINE TURN, PART ONE

## Does a language model know anything — or just bullshit?

[ESTABLISHED] [CONTESTED]

Day 1 ended on a needle of a question: when a system like the one that drafted these pages outputs a true, well-supported sentence, does it *know* anything — or is it the ultimate Gettier case, right for reasons that have nothing to do with the truth? The 2020s turned that closing flourish into one of the hottest debates in the field, and the most-discussed entry has a title that sailed past peer review unblunted: "ChatGPT is bullshit" (Hicks, Humphries & Slater, *Ethics and Information Technology*, 2024).

Their move is precise, not merely rude. They borrow Harry Frankfurt's technical sense of *bullshit* (from his 1986 essay *On Bullshit*): bullshit is speech produced with *indifference* to truth. The liar at least tracks the truth — he has to, in order to steer you away from it. The bullshitter doesn't care either way; he says whatever serves his purpose, and whether it's true is simply beside the point. Now consider what a large language model fundamentally *is*: a system trained to emit the statistically likely next token, to produce fluent, plausible-sounding text. It has no representation of truth that it is trying to honor. So when it states a real fact and when it "hallucinates" a fake citation, it is doing the *very same thing* — generating likely-looking text — and succeeding equally at its actual task in both cases. On this view, "hallucination" is a flattering misnomer that implies a malfunction; the truer description is that the system is **indifferent to truth by design**, which is bullshit in Frankfurt's exact sense. They distinguish *soft* bullshit (no intent to deceive, just truth-indifference) from *hard* bullshit (additionally posing as a sincere truth-teller), and argue an LLM is at minimum a soft bullshitter.

Why this could redraw the map: it cuts directly against the loose talk of machines "knowing," "understanding," or "believing." If the argument is right, an LLM's true outputs aren't knowledge and aren't even really *assertions* in the full sense — they're a new category of truth-apt-looking text with no one home who cares whether it's true. That reframes how

we should trust, cite, and regulate these systems. And it is, predictably, *contested* – the rebuttals are already a small literature. Some argue the "bullshit" label smuggles in a stance on whether models have intentions at all (Sarah Fisher, "Large language models and their big bullshit potential," 2024; David Gunkel & Simon Coghlan, "Cut the crap," 2025); others that as models are trained with reinforcement to be truthful and to express calibrated uncertainty, "indifferent to truth" is too crude. What's *settled* is the unglamorous core: a base language model has no built-in commitment to truth, and fluency is not knowledge. What's *open* is whether "bullshit," "instrument," "testifier," or some entirely new epistemic category is the right home for what these systems produce. It is the brain-in-a-vat from Appendix I made of silicon and shipped to a billion users – words that may never have touched the world, now answering our questions.

— §4 THE MACHINE TURN, PART TWO

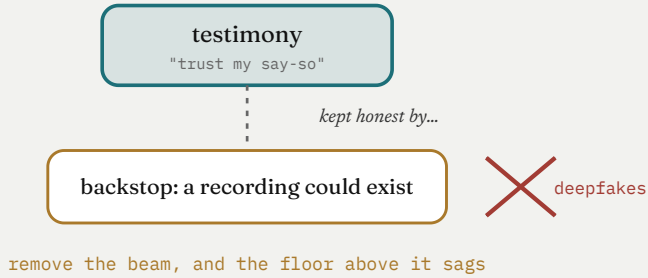
## The support beam you never noticed: the epistemic backstop

[PROMISING] [CONTESTED]

Appendix I established how much of what you know arrives by *testimony* – other people's say-so. Regina Rini's "Deepfakes and the Epistemic Backstop" (*Philosophers' Imprint*, 2020) identifies a hidden structural support that has quietly kept that whole edifice honest – and shows how a new technology is sawing through it.

The insight is subtle. Why is testimony as reliable as it is? Part of the answer, Rini argues, is a silent regulator: the ever-present *possibility* of a recording. When someone can be contradicted by a photograph, an audio clip, or a video, they have a standing incentive to testify truthfully – because a recording might surface and catch them out. Recordings function as an *epistemic backstop*: not because we constantly check them, but because their mere availability disciplines testimony, the way an unused referee still shapes a game. For roughly a century – since photography and audio became hard to fake – we have built our norms of public truth-telling on this backstop without ever naming it. Deepfakes – convincingly fabricated video and audio, generated by the same machine-learning wave as §3 – dissolve it from both sides. They flood the channel with convincing fakes *and*, just as corrosively, they hand every caught wrongdoer a new escape: *that recording of me is probably a deepfake*. The "liar's dividend." Once any recording can be waved away, the backstop stops disciplining testimony – and testimony itself, our single largest source of knowledge, loses a support we didn't know it leaned on. Don Fallis sharpens the same worry in information-theoretic terms ("The Epistemic Threat of Deepfakes," *Philosophy & Technology*, 2021):

deepfakes *reduce the information* a video carries about what actually happened, degrading it as a signal.



Knock out a support no one was looking at, and the structure it held still falls.

This is impactful precisely because it reframes deepfakes as an *epistemological* problem, not just a fraud or privacy one: the threat isn't only the specific lies, it's the erosion of a background condition for trusting recordings at all. But – true to this appendix – the magnitude is contested, and the pushback is sharp and worth taking seriously. Joshua Habgood-Coote ("Deepfakes and the epistemic apocalypse," *Synthese*, 2023) argues the doom framing is overblown: we have never relied on recordings as infallible, we already cross-check testimony against many sources, and societies have absorbed media-manipulation panics before. Atencia-Linares and Artiga ("Deepfakes, shallow epistemic graves," *Synthese*, 2022) defend the residual epistemic robustness of photography and video. So the *mechanism* Rini names – recordings as a silent regulator of testimony – is a genuine and illuminating contribution; the *prediction* of an "epistemic apocalypse" or wholesale collapse of public knowledge is a live dispute, not a settled forecast. Bring this one forward to **Day 2**, where science's answer to "how do you trust a report you can't personally verify?" is a whole institutional machinery of replication and recording – and to the AI block, where it meets §3 head-on.

## — §5 THE ADVERSARIAL TURN

# Hostile epistemology: when the environment is built to fool you

[ESTABLISHED] [CONTESTED]

Traditional epistemology pictured a lone, neutral mind facing a neutral world. C. Thi Nguyen's program – he calls it *hostile epistemology* – starts from a darker and more modern premise: your epistemic environment is not neutral. It is increasingly *engineered*, often by parties with an interest in what you come to believe, to exploit the predictable shortcuts your mind must use. Three of his post-2020 moves have reshaped how the whole field talks about online life.

The first is a distinction that sounds academic and turns out to be the key to everything: the difference between an *epistemic bubble* and an *echo chamber* ("Echo Chambers and Epistemic Bubbles," *Episteme*, 2020). They are *not* the same thing, and conflating them is why so many well-meaning fixes fail. In a bubble, outside voices are merely *absent* – you simply haven't been exposed to them (think a filter that only ever shows you agreeable sources). In a chamber, outside voices are *present but actively discredited* – you've been trained to *distrust* them in advance ("the mainstream media lies," "experts are corrupt"). The consequence is stark and counterintuitive, and the interactive below lets you feel it: the obvious intervention – *expose people to the other side* – pops a bubble but can *strengthen* a chamber, because inside a chamber, encountering the enemy's argument is exactly what the chamber predicted, and so confirms it.

## Bubble vs. Chamber, as exposure outcomes

STRUCTURE	OUTSIDE VOICES	EXPOSURE	OUTCOME	LESSON
Epistemic bubble	Absent, not refuted.	New sources can connect and	puncture the bubble.	Exposure can work when the problem is missing information.
Echo chamber	Present but pre-discredited.	Exposure can reinforce	distrust because the chamber predicted hostile outsiders.	The obvious repair can backfire when distrust is built into the structure.

The second move names a vulnerability inside your own head. In "The Seductions of Clarity" (*Royal Institute of Philosophy Supplement*, 2021), Nguyen argues that the *feeling* of

clarity – that satisfying click when everything seems to fall into place – functions as a *thought-terminating heuristic*. We use the sense that a matter has become clear as a signal that we've inquired enough and can stop. Usually fine. But it means clarity can be *weaponized*: a manipulator who can manufacture an exaggerated sense of clarity – a tidy ideology that explains everything, a conspiracy theory where every fact slots satisfyingly into place – can get you to *halt your inquiry early*, before you notice the holes. Notice how this snaps together with §1: clarity is dangerous precisely because it *terminates the zetetic process*. The slickest, most "it all makes sense now" account is, for that very reason, the one to interrogate hardest. The third move completes the toolkit: in "Trust as an Unquestioning Attitude" (*Oxford Studies in Epistemology*, 2022), Nguyen analyzes trust itself as the stance of *not questioning* – of taking something as a settled background you build on without re-checking. Indispensable (you can't re-derive everything from scratch), and exactly therefore exploitable: capture what someone trusts unquestioningly, and you've captured where they'll never think to look.

The honest tag here is twofold. The *conceptual* contributions – bubble vs. chamber, clarity as inquiry-terminating, trust as unquestioning – have been rapidly and widely adopted because they're genuinely clarifying and action-guiding. But two cautions earn their chips. First, philosophers are already pushing on the construct itself (Carey & Ventham, "There is no fresh air: a problem with the concept of echo chambers," *Episteme*, 2025). Second – and this is a hype-filter point the course insists on – the *empirical* social-science picture of how prevalent real-world echo chambers actually are is genuinely mixed; several large studies find most people's media diets are more varied than the "sealed echo chamber" image suggests. So treat the *conceptual machinery* as a sharp and durable tool, and the *empirical scale* of the phenomenon as an open measurement question. The framework is the contribution; the size of the fire it describes is still being measured.

## — §6 THE FORMAL REFOUNDATION

# Re-deriving Day 1's Dutch book — from truth, not money

[ESTABLISHED] [CONTESTED]

On Day 1 we justified the laws of probability with a *bribe*. The Dutch book theorem showed that if your credences break the probability rules, a clever bookie can sell you a set of bets you each accept as fair but which together guarantee you lose money. Powerful – but faintly unsatisfying as *epistemology*. Who cares about money? Shouldn't a *belief* be irrational for some reason to do with *truth*, not with your wallet? A research program that matured

through the 2010s and is in full bloom now – *accuracy-first epistemology*, also called epistemic utility theory – answers exactly that, and it's one of the most elegant results in the modern subject.

The idea (seeded by James Joyce's 1998 "Nonpragmatic Vindication of Probabilism" and built out by Richard Pettigrew's *Accuracy and the Laws of Credence*, 2016, with a wave of 2020–2023 papers refining and contesting it) is to measure how good a set of credences is by a single epistemic yardstick: *accuracy*, its closeness to the truth. Full confidence in a truth is perfectly accurate; full confidence in a falsehood, maximally *inaccurate*. Now the theorem. For any *incoherent* credence – one that violates the probability laws – there is guaranteed to exist a *coherent* credence that is **more accurate in every possible world at once**. The incoherent one is, in the technical term, *accuracy-dominated*: strictly beaten on truth-closeness no matter how things turn out. So you don't need the bookie at all. Incoherent confidence is irrational for a purely *epistemic* reason – it leaves guaranteed accuracy on the table; there's a better-aimed set of credences available that's closer to the truth whatever the world does. The interactive lets you see the domination geometrically.

## Accuracy domination, as credence geometry

CREDESCES	SUM	GEOMETRY	VERDICT
$P(S)=0.50, P(\text{not-}S)=0.50$	1.00	On the coherence line.	Undominated: no other credence is closer in every world.
$P(S)=0.80, P(\text{not-}S)=0.80$	1.60	Above the coherence line.	Dominated by a coherent projection closer to both truth-corners.
$P(S)=0.20, P(\text{not-}S)=0.20$	0.40	Below the coherence line.	Dominated by a coherent projection closer to both truth-corners.

What makes this a frontier rather than a footnote: it's an attempt to rebuild the *foundations* of rationality on a single epistemic value – getting close to the truth – and to derive not just probabilism but the update rule (conditionalization) and more besides from accuracy-dominance arguments. If it fully succeeds, the whole Bayesian edifice we started sketching

on Day 1 rests on truth, not on betting behavior or psychology. The chip, though, is earned. The *core* dominance theorem for probabilism is established mathematics. The *ambition* – that all epistemic norms fall out of accuracy alone – is contested: the cleanest results lean on technical assumptions (additivity, finitely many propositions) that critics argue smuggle in more than pure "closeness to truth" warrants (Chad Marxen, "Epistemic utility theory's difficult future," *Synthese*, 2021), and rival measures of accuracy can deliver different verdicts. So: a beautiful, genuinely illuminating reframing with a rock-solid center and a contested perimeter – which is, fittingly, the exact shape of this entire appendix.

## ◆ THE FRONTIER IN THREE SENTENCES

## BIG IDEA

Since 2020 the question "what is knowledge?" has been pushed from five directions at once — reconceiving epistemology as the study of *inquiry* not belief (zetetic), reordering the mind so knowledge comes *before* belief, and confronting machines, deepfakes, and engineered information environments that strain or attack the very notion of a knower — while a formal program quietly rebuilds rationality's foundations on *truth* itself.

## BEST NEW ANALOGY

The epistemic *backstop*: testimony has been kept honest all along by a support beam no one was looking at — the mere possibility of a recording — and deepfakes saw through it; pair it with the echo chamber, where the obvious fix (show them the other side) is precisely what makes the trap stronger.

## LIVE CONTROVERSY

Every item here is genuinely unsettled — whether inquiry-norms swallow belief-norms, whether knowledge-representation is really more basic than belief, whether "bullshit" is the right word for what LLMs do, whether deepfakes bring collapse or just friction, and whether accuracy alone can ground all of epistemic rationality — which is exactly why each wears a hype-filter tag.

---

THREADS HERE > information (testimony's hidden backstop; LLMs as truth-indifferent text engines; accuracy as the epistemic good) · computation (the mind's factive theory-of-mind; epistemic utility as decision theory for belief) · evolution (why a social species evolves to track *knowledge* first). The five threads, now at the waterline.

## — OPEN QUESTIONS

## What the edge of the map leaves blank

- **Is inquiry the real unit?** Do the norms of seeking truly conflict with the norms of believing – and if so, which is fundamental?
- **Knowledge or belief first?** Is "factive theory of mind" the basic cognitive tool, with belief a later, costlier add-on – or is the knowledge/belief line itself too crisp?
- **What *do machines produce*?** Knowledge, assertion, testimony, instrument-readings, or a genuinely new category of truth-apt text with no one home who cares?
- **Friction or collapse?** Do deepfakes merely add cost to verifying recordings, or dissolve a load-bearing condition of public knowledge?
- **How engineered is your mind's environment** – and how big, really, are the echo chambers we can now so clearly *describe*?
- **Can truth alone ground rationality?** Does accuracy-first reach all the way, or only as far as its technical assumptions carry it?
- **And a quieter contender for a future day:** several of these point past knowledge toward *understanding* as the thing we actually prize – a pivot we'll feel again whenever a model can predict without explaining.

## — SOURCES · ALL 2020+ UNLESS A FOUNDATIONAL ANCHOR

## Sources & further reading

1. Friedman, J. (2020). "The Epistemic and the Zetetic." *The Philosophical Review* 129(4): 501–536. doi:10.1215/00318108-8540918. link See also Falbo, A. (2023), "Should epistemology take the zetetic turn?" *Philosophical Studies* 180(10–11): 2977–3002; Flores, C. & Woodard, E. (2023), "Epistemic norms on evidence-gathering," *Philosophical Studies* 180(9): 2547–2571.
2. Phillips, J., Buckwalter, W., Cushman, F., Friedman, O., Martin, A., Turri, J., Santos, L. & Knobe, J. (2021). "Knowledge before belief." *Behavioral and Brain Sciences* 44: e140. doi:10.1017/S0140525X20000618 (target article + ~30 peer commentaries, several dissenting). link
3. Hicks, M. T., Humphries, J. & Slater, J. (2024). "ChatGPT is bullshit." *Ethics and Information Technology* 26: 38. doi:10.1007/s10676-024-09775-5. link Anchor: Frankfurt, H. (2005), *On Bullshit* (Princeton UP). Replies: Fisher, S. A. (2024), "Large language models and their big bullshit potential," *Ethics and*

- Information Technology* 26; Gunkel, D. & Coghlan, S. (2025), "Cut the crap: a critical response to 'ChatGPT is bullshit,'" *Ethics and Information Technology* 27.
4. Rini, R. (2020). "Deepfakes and the Epistemic Backstop." *Philosophers' Imprint* 20(24): 1–16. link And Fallis, D. (2021). "The Epistemic Threat of Deepfakes." *Philosophy & Technology* 34(4): 623–643. doi:10.1007/s13347-020-00419-2.
  5. Habgood-Coote, J. (2023). "Deepfakes and the epistemic apocalypse." *Synthese* 201(3). And Atencia-Linares, P. & Artiga, M. (2022). "Deepfakes, shallow epistemic graves: On the epistemic robustness of photography and videos in the era of deepfakes." *Synthese* 200(6). – the principal skeptical replies to the "collapse" framing.
  6. Nguyen, C. T. (2020). "Echo Chambers and Epistemic Bubbles." *Episteme* 17(2): 141–161. doi:10.1017/epi.2018.32. link
  7. Nguyen, C. T. (2021). "The Seductions of Clarity." *Royal Institute of Philosophy Supplement* 89: 227–255. And Nguyen, C. T. (2022). "Trust as an Unquestioning Attitude." *Oxford Studies in Epistemology* 7: 214–244. See also Nguyen (2023), "Hostile Epistemology," *Social Philosophy Today* 39: 9–32; and the critique Carey, B. & Ventham, E. (2025), "There is no fresh air: A problem with the concept of echo chambers," *Episteme* First View. doi:10.1017/epi.2024.43.
  8. Pettigrew, R. (2016). *Accuracy and the Laws of Credence*. Oxford University Press. Foundational anchor: Joyce, J. M. (1998), "A Nonpragmatic Vindication of Probabilism," *Philosophy of Science* 65(4): 575–603. Recent development & critique: Pettigrew, R. (2022), "Accuracy-First Epistemology Without Additivity," *Philosophy of Science* 89(1): 128–151; Marxen, C. (2021), "Epistemic utility theory's difficult future," *Synthese* 199(3–4): 7401–7421. Survey: SEP, "Epistemic Utility Arguments for Epistemic Norms."

Hype-filter note: classical anchors (Frankfurt 2005, Joyce 1998) are cited only as the roots of the post-2020 work that is this appendix's actual subject. No claim above should be treated as settled; that is the point of the chips.

TOMORROW → DAY 02

## The Scientific Method & Demarcation

Today we asked when a *single* belief counts as knowledge. Tomorrow we scale the question up to an entire institution: how does science decide which claims even get to enter the arena? Popper's demand that a real theory be *falsifiable*, Kuhn's paradigm shifts, Lakatos's research programmes – and the modern replication crisis as the demarcation line tested under live fire. Bring today's calibration instinct; you'll need it.

END OF DAY 01 · 179 DESCENTS REMAIN