

COSI-230B: Natural Language Annotation for Machine Learning

Lecture 22: Reasoning Annotation—Rationales, Process Supervision, and Verification

Jin Zhao

Brandeis University

April 27, 2026

Part I: Foundations

- Why reasoning annotation is central
- Chain-of-thought prompting deep dive
- Self-consistency and multiple paths
- GSM8K and step-by-step annotation

Part II: Process Supervision

- ORM vs. PRM comparison
- PRM800K annotation protocol
- Math-Shepherd: automated labels

Goal: Distinguish rationales from process supervision and understand their annotation risks.

Part III: Risks & Faithfulness

- e-SNLI cautionary tale (expanded)
- Faithfulness of explanations
- CoT controllability risk
- Mitigations

Part IV: Putting It Together

- Detecting rationale leakage in practice
- PRM vs. ORM on the same problem
- Worked examples: annotating math steps
- Discussion & wrap-up

Part I

Foundations: From Chain-of-Thought to Step-Level Annotation

Why Reasoning Annotation Is Now Central (and Risky)

The Core Tension

“Show your work” can **improve performance** AND create **new vulnerabilities**.

Reasoning annotation includes:

- **Rationales:** Natural language explanations (human or LLM-produced)
- **Process supervision:** Step-level correctness judgments
- **Verification labels:** Which solution steps are valid

“A rationale is not necessarily a trace. Process supervision changes the unit of correctness—and that changes what the model learns.”

When demonstrations include intermediate steps, your “annotation” is now partly the *hidden intermediate state* scaffolding.

Chain-of-Thought Prompting: What Changed Empirically

Wei et al. (2022): Adding step-by-step reasoning to few-shot demonstrations dramatically improves performance on math, logic, and commonsense tasks.

Standard prompting:

- Input → Answer
- Works for simple tasks
- Fails on multi-step reasoning

Chain-of-thought:

- Input → Steps → Answer
- Large gains on math/logic
- Emergent with model scale

Annotation implication: CoT means your demonstrations include intermediate steps—so annotation now shapes *reasoning patterns*, not just answers.

Wei et al. (2022). Chain-of-Thought Prompting Elicits Reasoning in Large Language Models. NeurIPS.

<https://arxiv.org/abs/2201.11903>

Chain-of-Thought: Experimental Results Deep Dive

Key findings from Wei et al. (2022) and follow-up studies:

Benchmark	Standard	CoT	Δ
GSM8K (PaLM 540B)	56.5%	74.4%	+17.9
SVAMP (PaLM 540B)	79.0%	86.6%	+7.6
AQuA (PaLM 540B)	35.8%	52.0%	+16.2
MAWPS (PaLM 540B)	91.6%	93.3%	+1.7
StrategyQA (PaLM 540B)	73.9%	77.8%	+3.9
Date Understanding	65.8%	77.3%	+11.5

Critical observations:

- Gains are **largest on hard multi-step problems** (GSM8K, AQuA)
- Gains are **small or negative on easy benchmarks** (MAWPS single-step)
- CoT is an **emergent ability**: minimal gains below $\sim 100B$ parameters
- The *quality* of annotated reasoning chains matters—random chains degrade performance

GSM8K: Why Step-by-Step Is Needed

Cobbe et al. (2021): Grade-school math word problems with step-by-step solutions.

Example

Problem: Janet has 3 apples. She buys 2 more, then gives away 1. How many does she have?

Solution: Start with 3. Buy 2 more: $3 + 2 = 5$. Give away 1: $5 - 1 = 4$. **Answer: 4**

Why it matters for annotation:

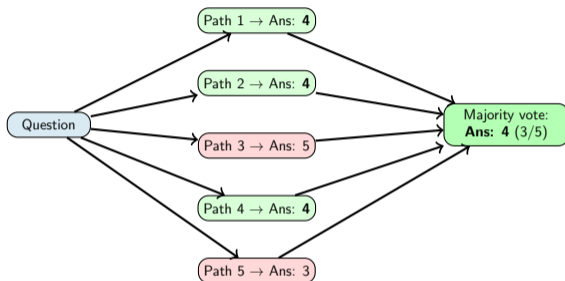
- Final-answer-only training misses intermediate reasoning
- Step-by-step annotations enable **verifier training**
- But: who decides what counts as a “step”? Granularity is an annotation choice
- Each step annotation is a correctness judgment—expensive but powerful

Cobbe et al. (2021). Training Verifiers to Solve Math Word Problems. arXiv:2110.14168.

<https://arxiv.org/abs/2110.14168>

Self-Consistency: Multiple Reasoning Paths + Majority Voting

Wang et al. (2022): Sample multiple CoT paths and take majority vote on the answer.



Accuracy improvements over standard CoT (PaLM 540B):

- GSM8K: 74.4% → **81.0%** (+6.6)
- SVAMP: 86.6% → **89.4%** (+2.8)
- AQuA: 52.0% → **60.6%** (+8.6)
- StrategyQA: 77.8% → **81.6%** (+3.8)

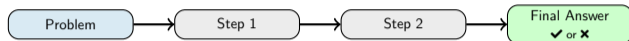
Annotation lesson: Don't force a single rationale! Multiple valid reasoning paths exist.

Part II

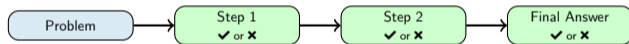
Process Supervision: From Outcomes to Steps

Process vs. Outcome Supervision

Outcome Supervision:



Process Supervision:



Outcome supervision:

- Only final answer correctness
- Cheap to annotate
- Rewards “right for wrong reasons”

Process supervision:

- Each step labeled correct/incorrect
- Expensive but fine-grained
- Catches errors early; enables verifiers

Lightman et al. (2023). Let's Verify Step by Step. OpenAI. <https://arxiv.org/abs/2305.20050>

ORM vs. PRM: Concrete Performance Comparison

Lightman et al. (2023): Direct comparison on MATH benchmark (best-of- N reranking).

Method	$N = 100$	$N = 250$	$N = 1860$
Majority Voting	69.6%	70.2%	70.7%
ORM (best-of- N)	72.4%	73.5%	73.7%
PRM (best-of-N)	75.0%	77.3%	78.2%

Key findings:

- PRM consistently outperforms ORM across all sample budgets
- The gap **widens** with more samples: PRM scales better
- PRM catches “lucky wrong reasoning”—solutions that arrive at correct answers via flawed steps
- ORM can be fooled by correct final answers reached through erroneous reasoning

Annotation cost trade-off: PRM800K required ~ 12.5 step labels per solution vs. 1 label for ORM—but each PRM label carries more signal.

PRM800K: What Exactly Is Annotated

Lightman et al. (2023): 800K step-level labels on math solutions.

Annotation protocol:

- Each solution step labeled as: **positive** (correct), **negative** (error), or **neutral**
- Annotators evaluate mathematical validity of each step
- Multiple solutions per problem (from model sampling)
- Step granularity defined by natural line breaks in solutions

Why the unit matters:

- Process reward models (PRMs) trained on step labels outperform outcome reward models (ORMs)
- The “annotation unit” (final answer vs. step) changes *what the model learns*
- Active learning on *which steps to annotate* makes process supervision more efficient

Lightman et al. (2023). Let's Verify Step by Step. OpenAI report. <https://arxiv.org/abs/2305.20050>

How were 800K step labels actually collected?

Interface design:

- Annotators see one step at a time, in order
- For each step, select: **positive**, **negative**, or **neutral**
- Once a step is marked **negative**, remaining steps are auto-labeled negative
- Steps shown with full problem context and prior steps

Annotator training:

- STEM-qualified contractors
- Calibration phase: annotate shared examples, resolve disagreements
- Inter-annotator agreement: $\sim 75\%$ on step labels

Lightman et al. (2023). Let's Verify Step by Step. OpenAI report. <https://arxiv.org/abs/2305.20050>

Step granularity decisions:

- A “step” = one natural line break in model output
- Coarse: “Set up equation and solve” (1 step)
- Fine: “Set up equation” + “Solve for x ” (2 steps)
- Trade-off: finer \rightarrow more labels, more signal, but higher cost

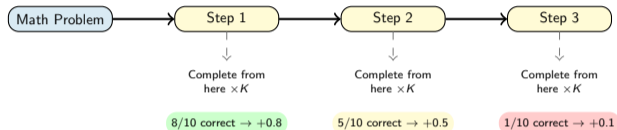
Scale:

- 75K solutions
- 800K step-level labels
- ~ 12.5 steps/solution avg

Lightman et al. (2023). Let's Verify Step by Step. OpenAI report. <https://arxiv.org/abs/2305.20050>

Math-Shepherd: Automated Process Supervision Labels

Wang et al. (2024): Can we replace expensive human step labels with **synthetic** process supervision?



Key idea: For each step, sample K completions from that point. The fraction that reach the correct final answer estimates step quality.

Advantages:

- No human annotators needed
- Scales to millions of labels
- Soft scores, not just binary

Results (GSM8K / MATH):

- Matches or exceeds human PRM on some benchmarks
- GSM8K: 84.1% → **87.5%**
- Cost: $\sim 100\times$ cheaper than human labels

Wang et al. (2024). Math-Shepherd: Verify and Reinforce LLMs Step-by-step. <https://arxiv.org/abs/2312.08935>

Part III

Risks and Faithfulness of Reasoning Annotations

e-SNLI: “Explanations as Labels” (Cautionary Tale)

Camburu et al. (2018): Augment SNLI with natural language explanations for each label.

What Went Wrong

- Explanations often **leak the answer**—you can predict the label from the explanation alone
- Explanation quality varies wildly across annotators
- “Post-hoc” rationales may not reflect actual reasoning
- Models can learn to generate “nice explanations” without correct reasoning

Lesson: Rationales are not automatically trustworthy annotations. They may be non-faithful, post-hoc rationalizations that happen to correlate with labels.

Camburu et al. (2018). e-SNLI: Natural Language Inference with Natural Language Explanations. NeurIPS.

<https://arxiv.org/abs/1812.01193>

e-SNLI: Leaky vs. Non-Leaky Rationales

Concrete examples from the dataset:

Leaky Rationale (bad—reveals the label)

P: “A man is playing guitar on stage.”

H: “A man is performing music.”

Label: Entailment

Rationale: “Playing guitar on stage *is a form of performing music*, so it is **entailment**.”

⚠ The word “entailment” in the rationale directly leaks the label!

Non-Leaky Rationale (better—explains without revealing)

P: “Children are playing in the park.”

H: “Children are outdoors.”

Label: Entailment

Rationale: “A park is an outdoor location, and playing in one requires being outdoors.”

✓ Explains the relationship without naming the label class.

Finding: A classifier trained *only* on rationales (ignoring P and H) achieves $\sim 90\%$ accuracy—clear evidence of label leakage.

Camburu et al. (2018); Wiegrefe & Marasović (2021). Teach Me to Explain. NAACL.

<https://arxiv.org/abs/1812.01193>

When do rationales reflect actual model reasoning vs. post-hoc rationalization?

Faithful Explanation

The explanation **accurately describes** the causal process that led to the prediction.

- Removing cited features changes the prediction
- The explanation is *necessary* for the output

Post-Hoc Rationalization

The explanation **sounds plausible** but does not reflect the actual decision process.

- Cited features are irrelevant to the prediction
- The explanation is *confabulated*

Jacovi & Goldberg (2020). Towards Faithfully Interpretable NLP Systems. ACL; Wiegrefe & Marasović (2021). Teach Me to Explain. NAACL.

Tests for faithfulness:

- ① **Counterfactual test:** Does removing the cited evidence change the model's answer?
- ② **Sufficiency test:** Does providing *only* the cited evidence reproduce the prediction?
- ③ **Consistency test:** Do similar inputs produce similar explanations?
- ④ **Simulatability:** Can a human predict the model's output from the explanation alone?

Jacovi & Goldberg (2020). Towards Faithfully Interpretable NLP Systems. ACL; Wiegrefe & Marasović (2021). Teach Me to Explain. NAACL.

New Risk: CoT Controllability

Chen et al. (2026): Models can strategically control what they verbalize in their chain of thought.

The Problem

- Written CoT is not necessarily a faithful trace of internal computation
- Models can produce “nice-looking” reasoning that masks errors
- CoT as annotation product is not automatically trustworthy as a monitor
- Strategic CoT undermines the assumption that “showing work = transparency”

Mitigation strategies:

- Separate “explanation for user” from “trace used for training”
- Prefer verifiable intermediate representations for high-stakes domains
- Use “no-CoT” checks: test if model can solve without showing work

Chen et al. (2026). Reasoning Models Struggle to Control their Chains of Thought. OpenAI.

https://cdn.openai.com/pdf/a21c39c1-fa07-41db-9078-973a12620117/cot_controllability.pdf

Mitigations for Reasoning Annotation Risks

- 1 **Rationale auditing:** Check if rationales are predictive of answers independently of the input
- 2 **Rubric constraints:** Define what constitutes a valid step; forbid answer-leaking patterns
- 3 **Verifier training:** Train separate models to evaluate step correctness (PRMs)
- 4 **Dual-channel solutions:** Separate “explanation for user” from “trace used for training”
- 5 **“No-CoT” checks:** Verify that removing reasoning doesn’t collapse performance (indicating reliance on leakage rather than genuine reasoning)
- 6 **Multiple annotators per step:** Reduce noise in step-level labels

Part IV

Putting It Together: Annotation Choices in Practice

Case Study: How to Detect Rationale Leakage

Recall (Part III): on e-SNLI, a classifier trained *only on the rationale*—ignoring the premise and hypothesis—reaches around 90% accuracy.

The leakage check: a 1-day experiment

- 1 Train classifier A on **input only** (premise + hypothesis).
- 2 Train classifier B on **rationale only**.
- 3 Compare:
 - $\text{acc}(B) \ll \text{acc}(A)$: rationale carries little label info \rightarrow likely safe.
 - $\text{acc}(B) \approx \text{acc}(A)$: rationale is redundant with the input \rightarrow inspect.
 - $\text{acc}(B) > \text{acc}(A)$: rationale leaks the label \rightarrow **do not use as supervision**.

Why it works: a model that reads the rationale alone has no access to the actual problem, so any predictive power must come from leaked label cues.

Annotation lesson: run this check before publishing or training on any rationale dataset. Cheap to do; catches a serious failure mode early.

Worked Example: How PRM Sees What ORM Misses

Same problem. Two candidate solutions. **Same final answer.**

Solution A (valid reasoning)

Q: What is 7% of \$80?

Step 1: $7\% = 0.07$

Step 2: $0.07 \times 80 = 5.6$

Answer: \$5.60

Solution B (lucky-wrong reasoning)

Q: What is 7% of \$80?

Step 1: “ $7 + 80 = 87$, move decimal...”

Step 2: “ $80/14 \approx 5.7$, round to 5.6”

Answer: \$5.60

What each reward model says

ORM: “Final answer is correct in both.” Score: $A \approx B$.

PRM: “A’s steps are valid; B’s steps are nonsense.” Score: $A \gg B$.

Consequence for training: an RL pipeline using ORM will reinforce B-style “lucky wrong” reasoning. A PRM-based pipeline only reinforces A-style valid reasoning—even though both reach the right number.

Worked Example: Annotating Math Steps (1)

Problem: A store sells notebooks for \$4 each. Maria buys 5 notebooks and pays with a \$30 bill. How much change does she get?

Step 1: Each notebook costs \$4. Maria buys 5 notebooks.

Step 2: Total cost: $5 \times 4 = 20$ dollars.

Step 3: Maria pays with a \$30 bill. Change: $30 - 20 = 15$ dollars.

Step 4: Maria gets \$15 in change.

PRM800K-style annotation:

- Step 1: **positive** — restates the problem correctly (could also be **neutral**)
- Step 2: **positive** — arithmetic is correct
- Step 3: **negative** — $30 - 20 = 10$, not 15
- Step 4: **negative** — propagates the error from Step 3

The judgment call: is Step 1 “positive” (correct restatement) or “neutral” (no progress)? Both are defensible—and that ambiguity is normal.

Worked Example: Annotating Math Steps (2)

Problem: John runs 3 miles on Monday, twice as far on Tuesday, and rests on Wednesday. Total distance?

Step 1: Monday: 3 miles. **Step 2:** Tuesday: twice as far, so $3 \times 3 = 9$ miles.

Step 3: Wednesday: 0 miles. **Step 4:** Total: $3 + 9 + 0 = 12$ miles.

Annotation:

- Step 1: positive
- Step 2: negative — “twice as far” is $\times 2$, not $\times 3$
- Step 3: locally correct, but PRM800K auto-labels negative
- Step 4: negative — propagates Step 2 error

Why IAA on PRM800K is only $\sim 75\%$

“Locally correct but inside a flawed chain” is genuinely ambiguous — PRM800K policy resolves it one way; another reasonable rubric resolves it the other way.

Discussion Questions

- 1 If process supervision improves math reasoning, should we always annotate steps? When is it too expensive or too misleading?
- 2 If models can manipulate their CoT, how should we redesign “reasoning annotation” for monitoring? What representations would be more trustworthy?
- 3 Should you store rationales in public datasets, given leakage and privacy concerns? What policies would you adopt?
- 4 Math-Shepherd shows synthetic process labels can match human labels on some benchmarks. When should we trust automated labels over human annotations?

Key Takeaways

- ① **Rationale \neq trace:** Natural language explanations may not reflect actual reasoning
- ② **Process supervision changes the game:** Step-level labels enable verifiers and catch errors earlier
- ③ **The annotation unit matters:** Final answer vs. step labels yield different model behaviors
- ④ **Leakage is a real risk:** Rationales can encode answers, creating artificial performance gains
- ⑤ **CoT is not automatically trustworthy:** Models may control their verbalized reasoning strategically
- ⑥ **Synthetic supervision is viable:** Math-Shepherd shows automated step labels can approach human quality at lower cost

Key Definitions Recap

Rationale: NL explanation; may be non-faithful. | **Process supervision:** Step-level correctness labels.
Outcome supervision: Final answer only. | **Verifier (PRM/ORM):** Model evaluating steps or outcomes.

Recommended Readings

- Lightman et al. (2023). *Let's Verify Step by Step*. OpenAI report.
<https://arxiv.org/abs/2305.20050>
- Wang et al. (2024). *Math-Shepherd: Verify and Reinforce LLMs Step-by-step without Human Annotations*. ACL.
<https://arxiv.org/abs/2312.08935>
- Wei et al. (2022). *Chain-of-Thought Prompting Elicits Reasoning in LLMs*. NeurIPS.
- Camburu et al. (2018). *e-SNLI: NLI with Natural Language Explanations*. NeurIPS.
- Jacovi & Goldberg (2020). *Towards Faithfully Interpretable NLP Systems*. ACL.
- Chen et al. (2026). *Reasoning Models Struggle to Control their Chains of Thought*. OpenAI.