

Distill or Annotate? Cost-Efficient Fine-Tuning of Compact Models



Junmo Kang



Wei Xu



Alan Ritter

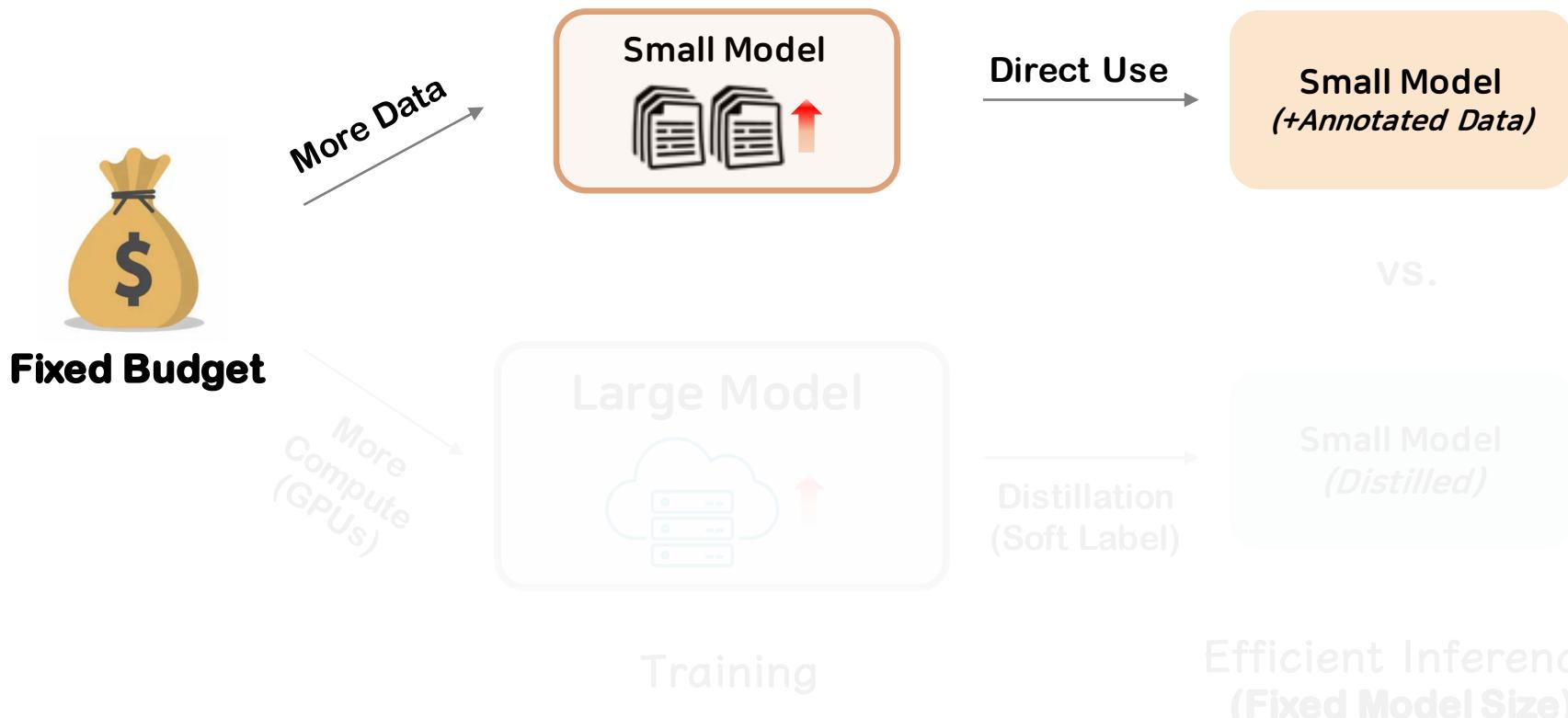


*Q. Given a fixed budget,
how to build a compact model in a cost-efficient way?*

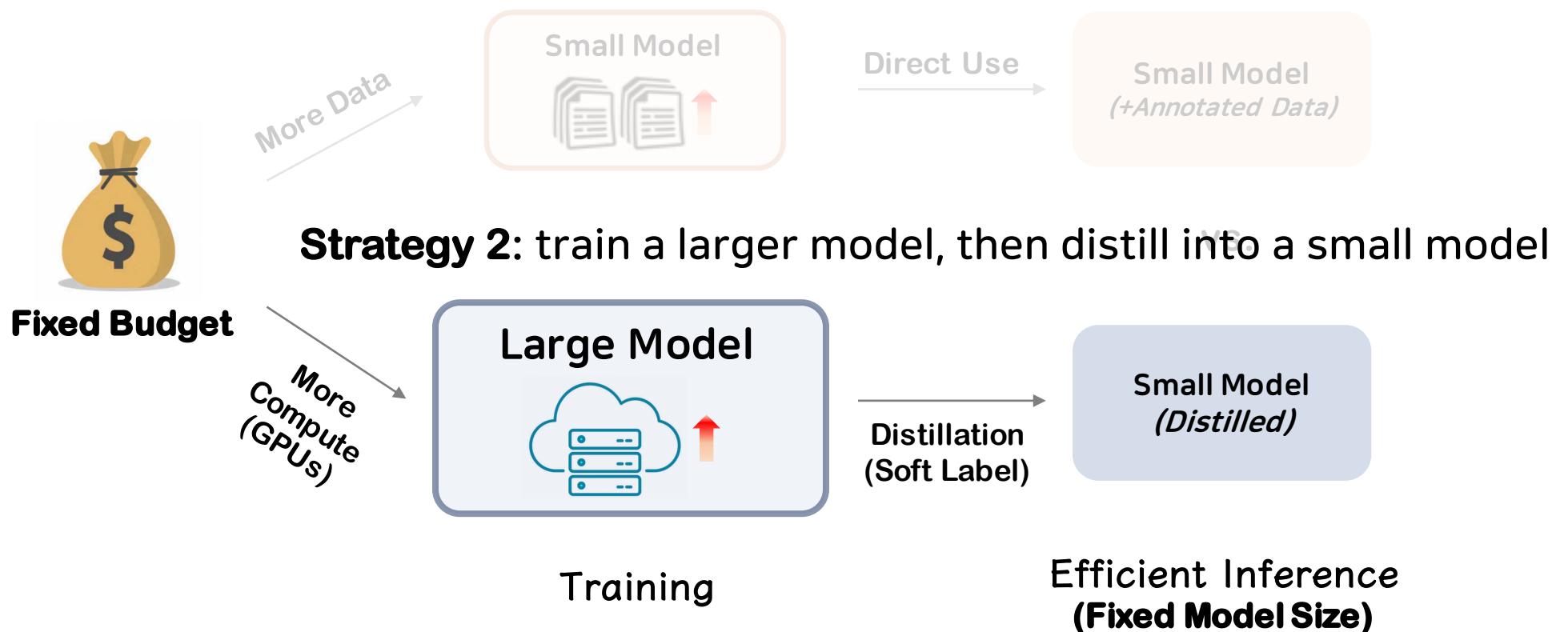


*Q. Given a fixed budget,
how to build a compact model in a cost-efficient way?*

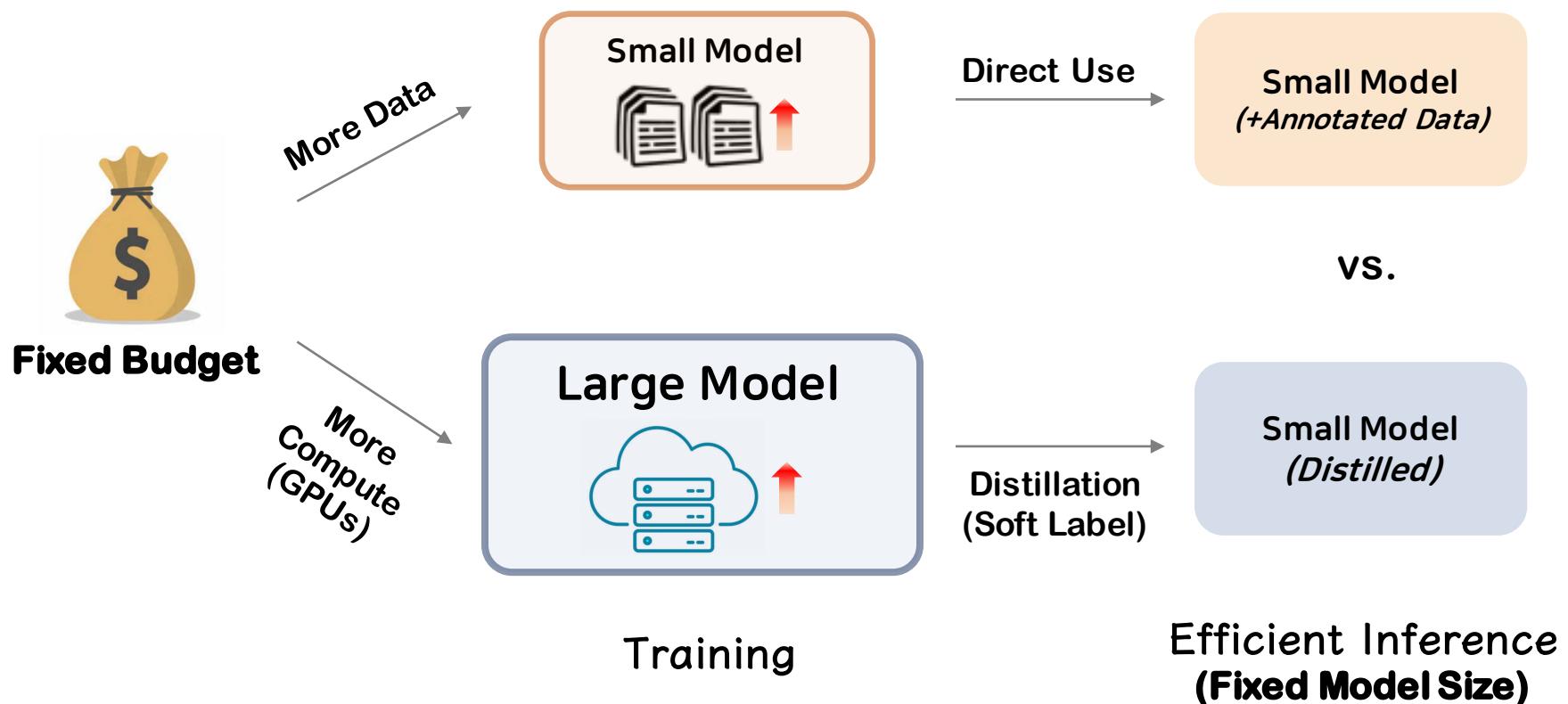
Strategy 1: annotate more data to directly train a small model



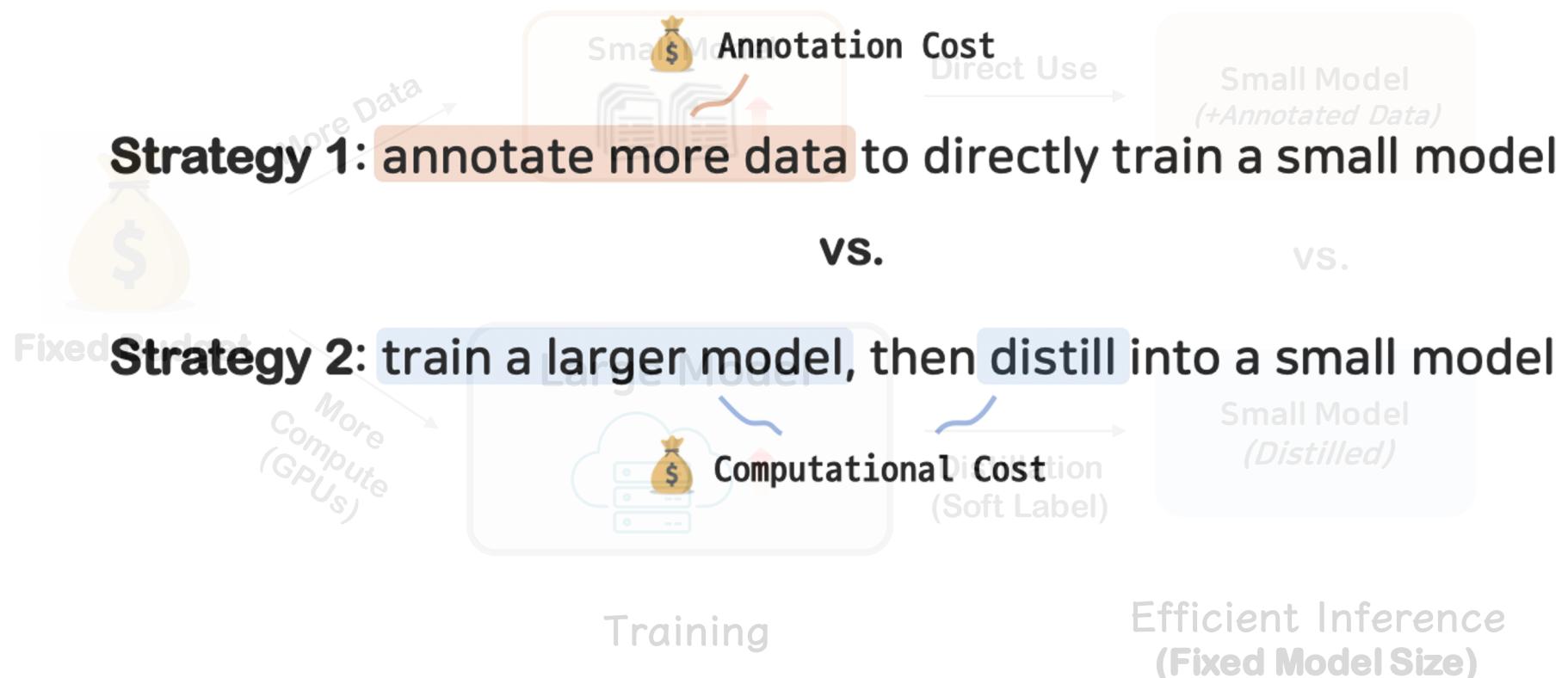
*Q. Given a fixed budget,
how to build a compact model in a cost-efficient way?*



*Q. Given a fixed budget,
how to build a compact model in a cost-efficient way?*



Trade-Off



Cost Estimations

Task & Annotation Cost

Dataset	Task	\$ per Label
WLP	Named Entity Recognition	\$0.26
STANCEOSAURUS	Stance Classification	\$0.364
FEVER	Fact Verification	\$0.129
MULTIPIT _{ID}	Paraphrase Identification	\$0.2
MULTIPIT _{GEN}	Paraphrase Generation	\$0.371
Natural Questions	Question Answering	\$0.129

Cost Estimations

Task & Annotation Cost

Dataset	Task	\$ per Label
WLP	Named Entity Recognition	\$0.26
STANCEOSAURUS	Stance Classification	\$0.364
FEVER	Fact Verification	\$0.129
MULTIPIT _{ID}	Paraphrase Identification	\$0.2
MULTIPIT _{GEN}	Paraphrase Generation	\$0.371
Natural Questions	Question Answering	\$0.129

Computational Cost

\$1.875 per 1 GPU hour (est. based on A100 in Google Cloud Platform)



Main Results

WLP (F1)

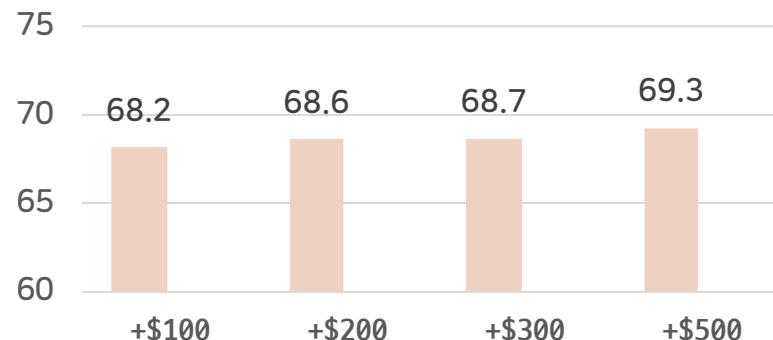
Starting #Data: 5K (\$1300)



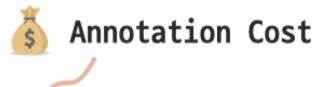
Main Results

WLP (F1)

Starting #Data: 5K (\$1300)

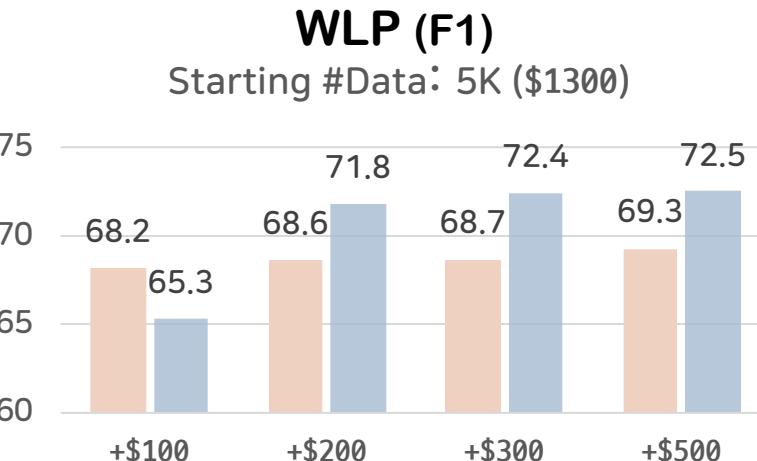


■ T5-Small (Ann.)



Strategy 1: **annotate more data** to directly train a small model

Main Results



- T5-Small (Ann.)
- T5-XXL => T5-Small (Dist.)

💰 Annotation Cost

Strategy 1: **annotate more data** to directly train a small model

vs.

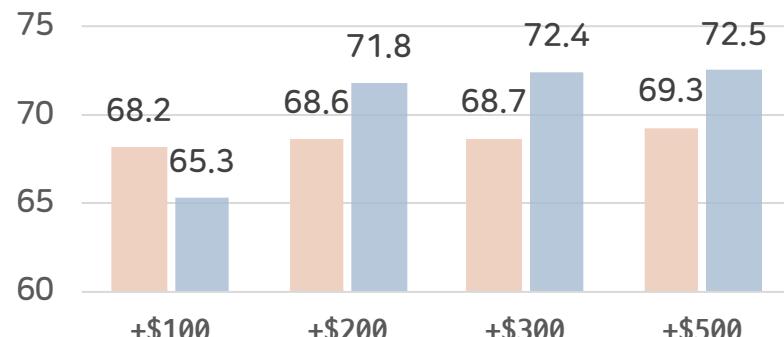
💰 Computational Cost

Strategy 2: **train a larger model**, then **distill** into a small model

Evaluation



WLP (F1)
Starting #Data: 5K (\$1300)



MULTIPIT_{Id} (Accuracy)
Starting #Data: 5K (\$1000)



Main Results

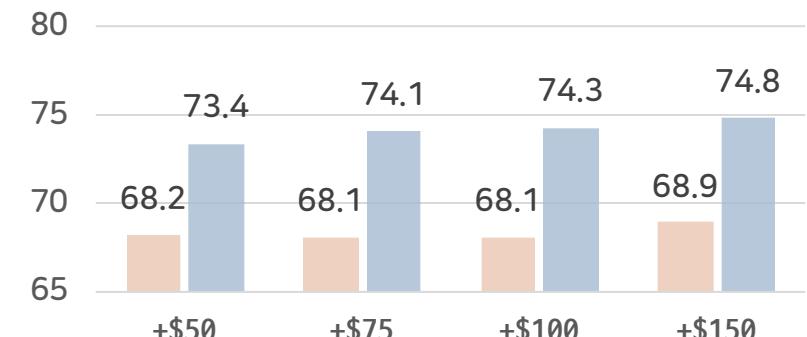
STANCEOSAURUS (F1)
Starting #Data: 5K (\$1820)



MULTIPIT_{Gen} (BERT-iBLEU)
Starting #Data: 10K (\$3710)



FEVER (Accuracy)
Starting #Data: 5K (\$645)



NATURAL QUESTIONS (F1)
Starting #Data: 10K (\$1290)



■ T5-Small (Ann.) ■ T5-XXL => T5-Small (Dist.)

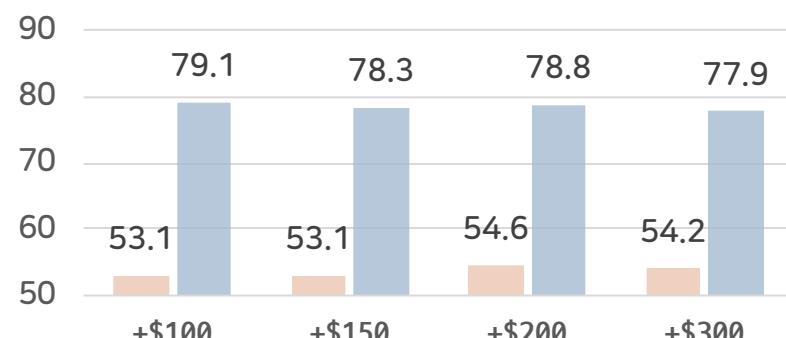
WLP (F1)

Starting #Data: 1K (\$260)



MULTIPIT_{Id} (Accuracy)

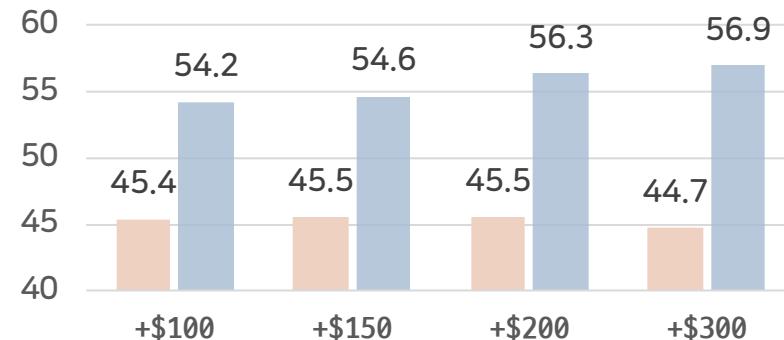
Starting #Data: 1K (\$200)



Main Results

STANCEOSAURUS (F1)

Starting #Data: 1K (\$364)



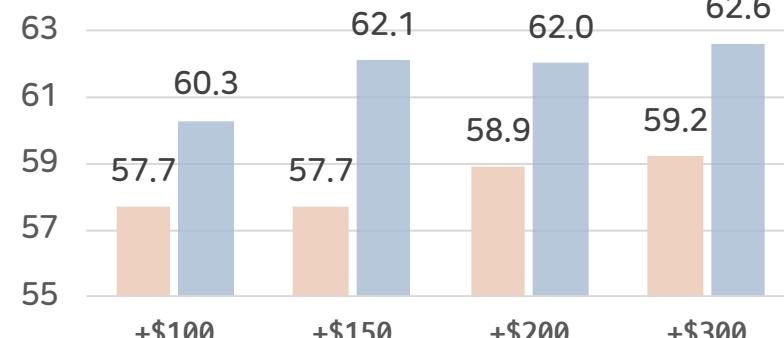
FEVER (Accuracy)

Starting #Data: 1K (\$129)



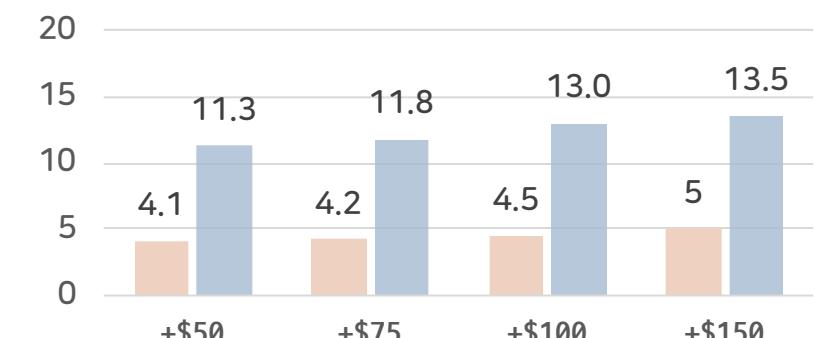
MULTIPIT_{Gen} (BERT-iBLEU)

Starting #Data: 1K (\$371)



NATURAL QUESTIONS (F1)

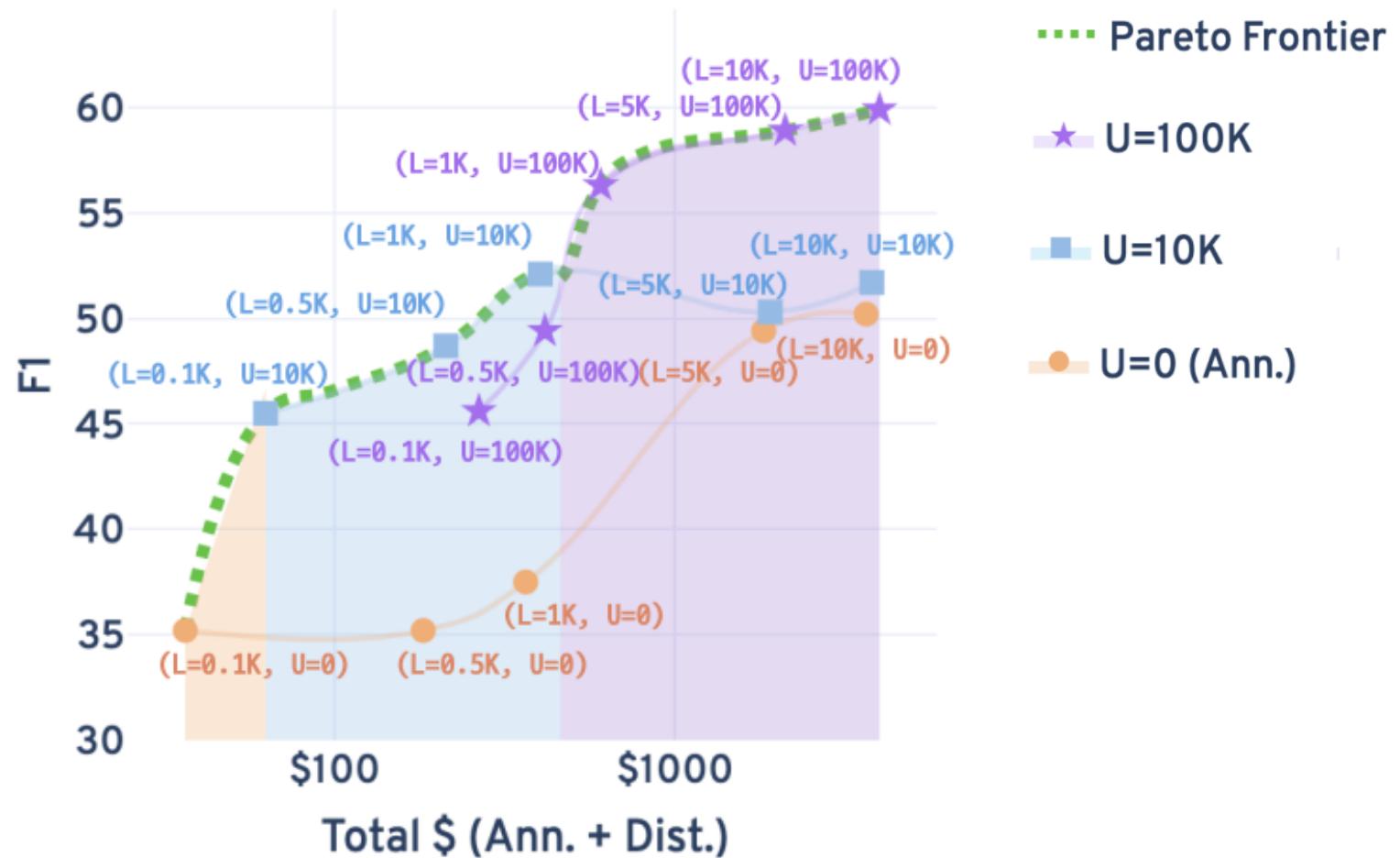
Starting #Data: 1K (\$129)



■ T5-Small (Ann.) ■ T5-XXL => T5-Small (Dist.)

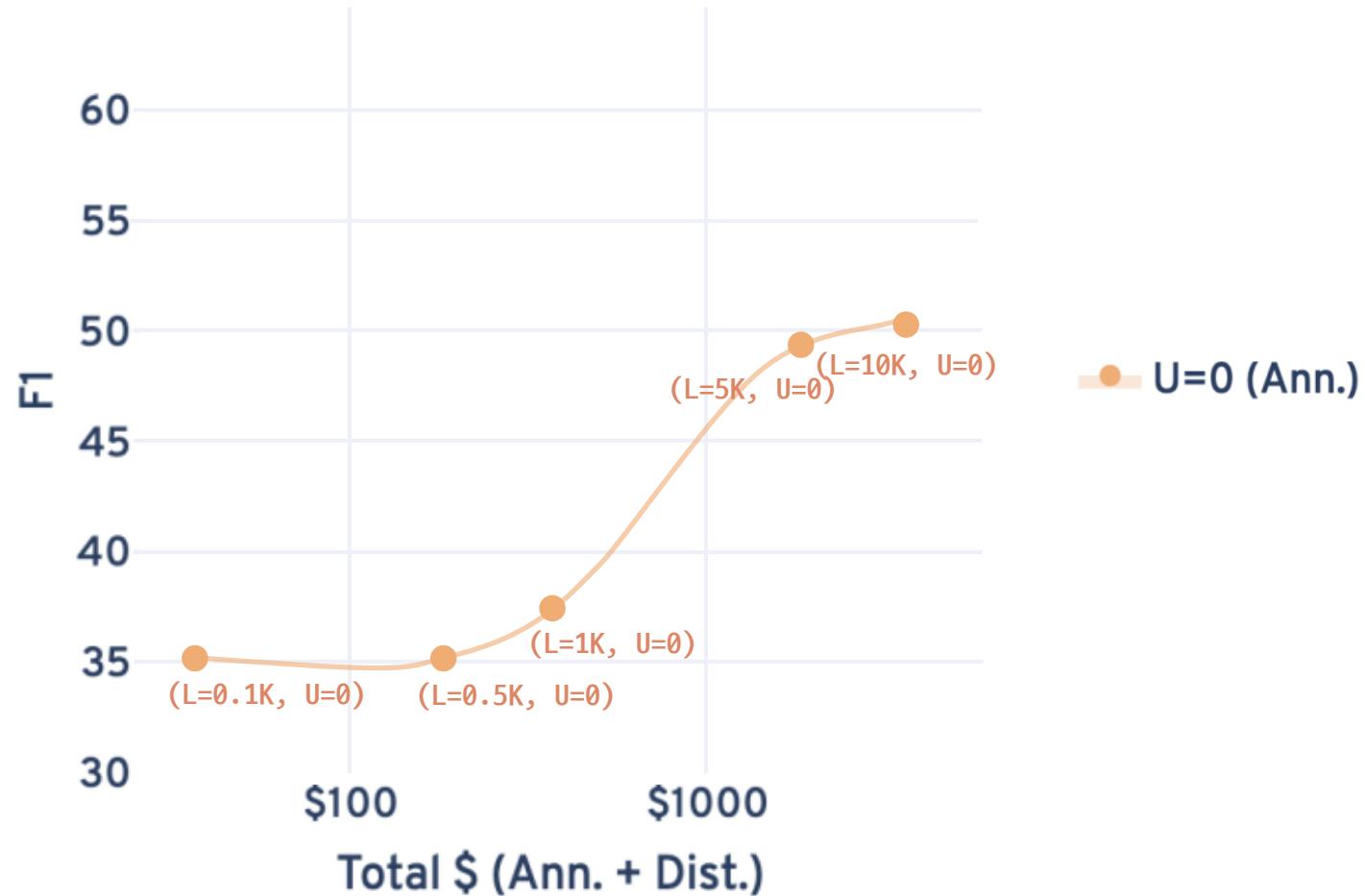
Pareto Curves

STANCEOSAURUS



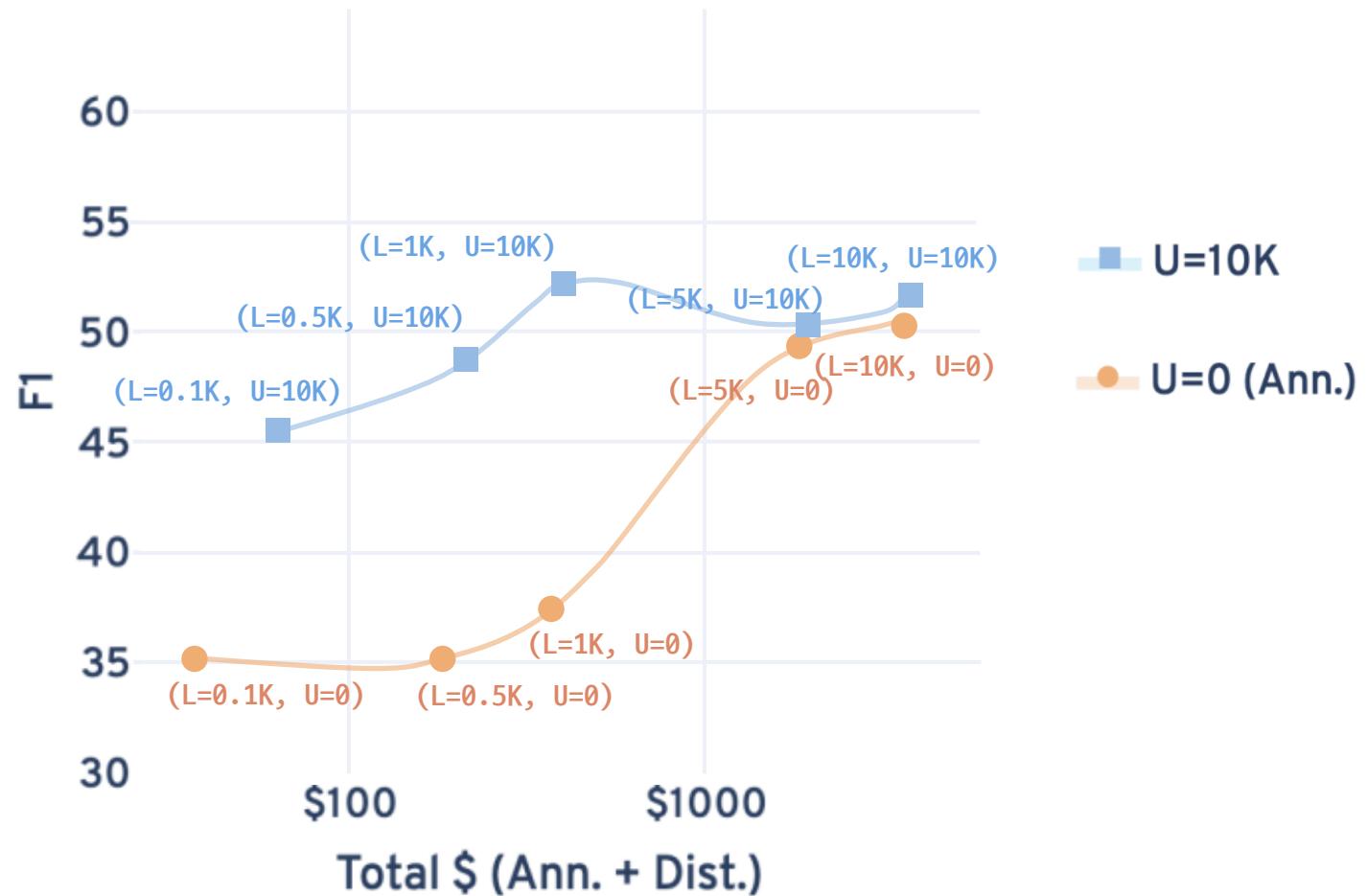
Pareto Curves

STANCEOSAURUS



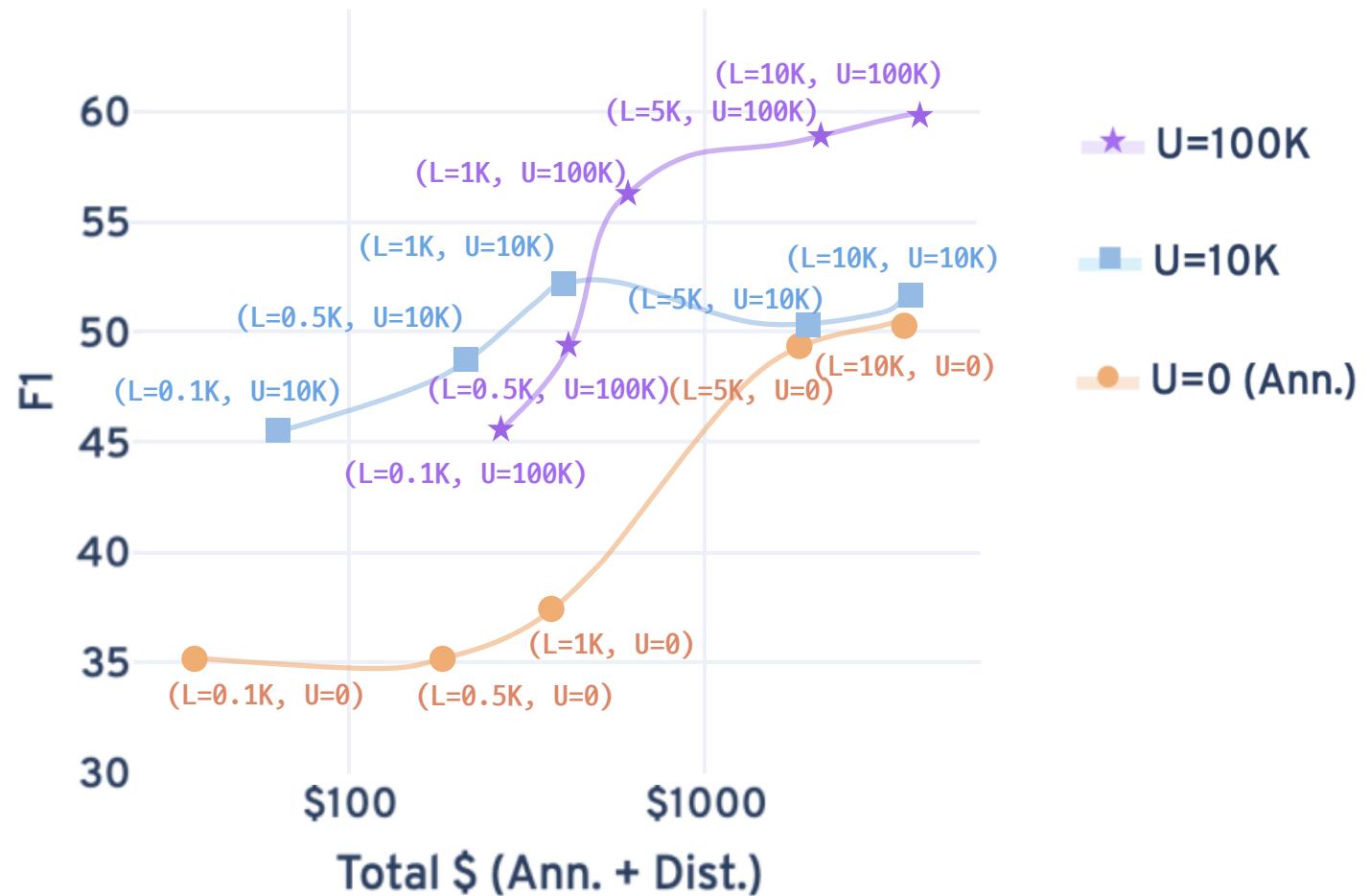
Pareto Curves

STANCEOSAURUS



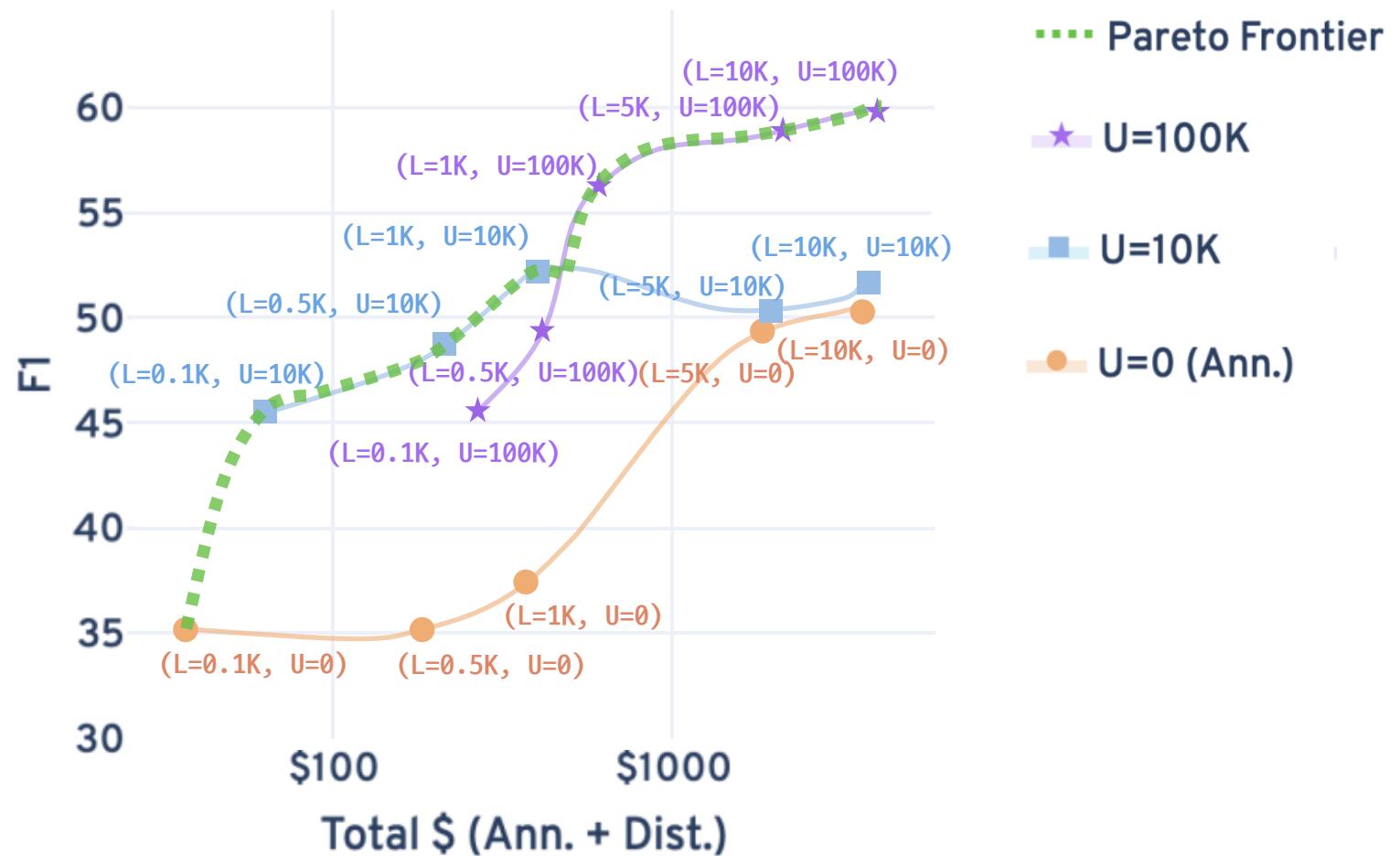
Pareto Curves

STANCEOSAURUS



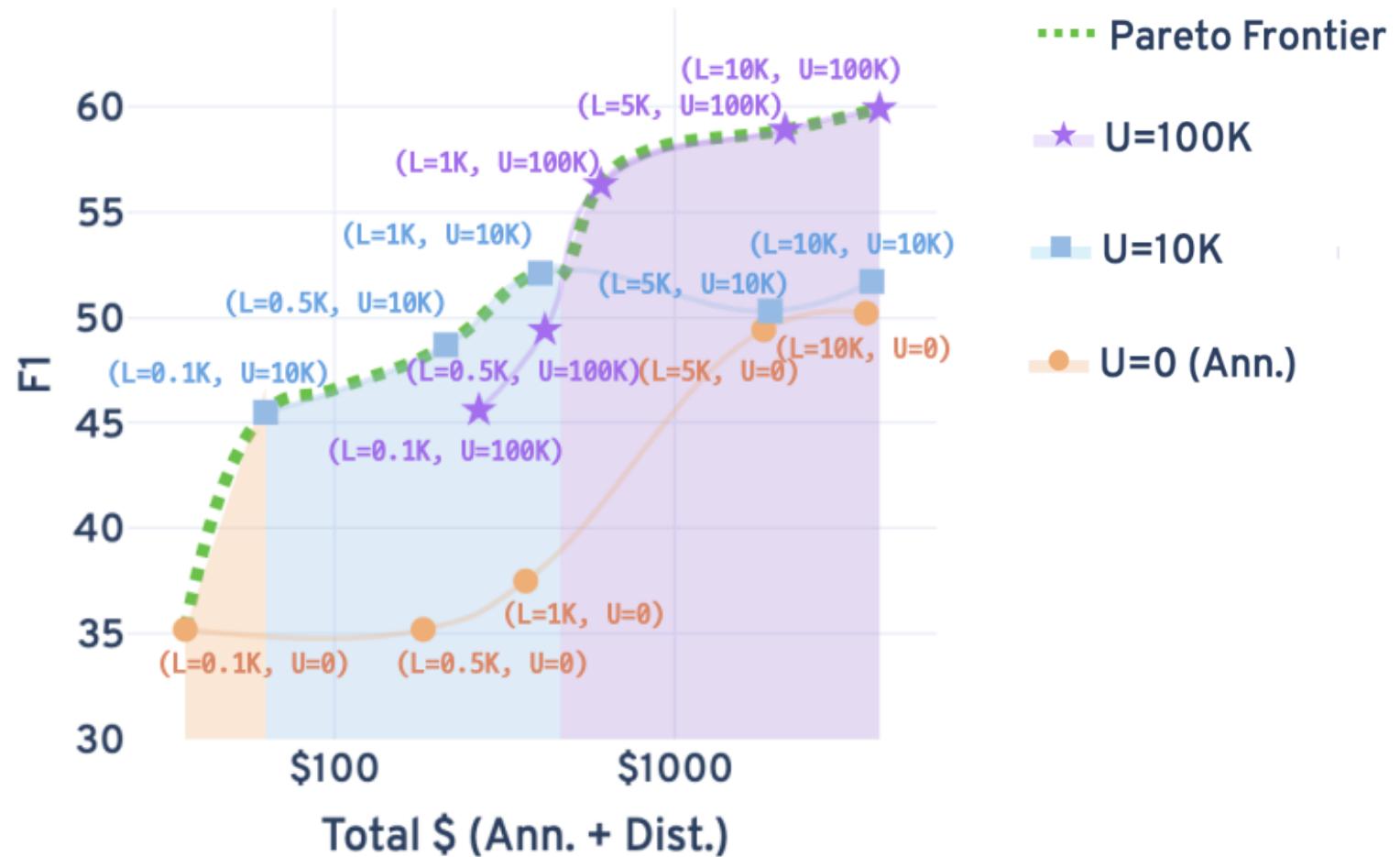
Pareto Curves

STANCEOSAURUS



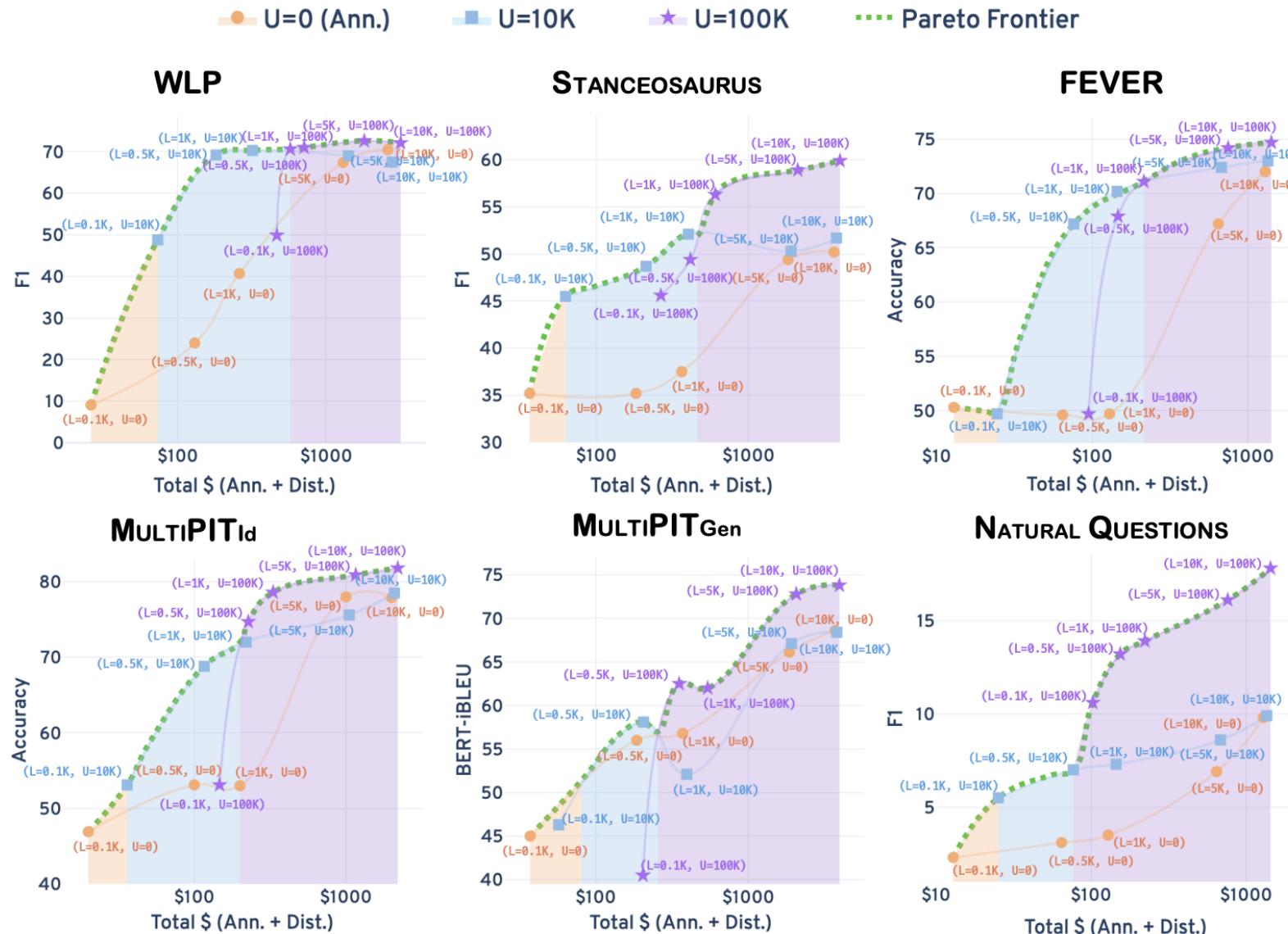
Pareto Curves

STANCEOSAURUS



Evaluation

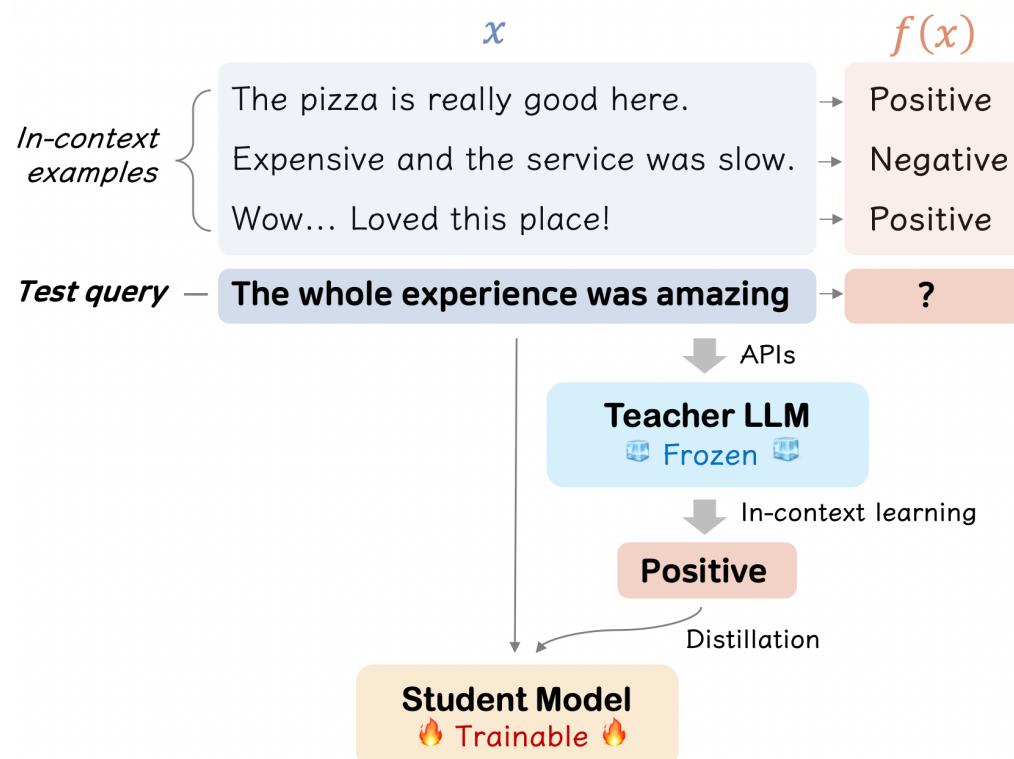
Pareto Curves



GPT-3.5 as an Annotator

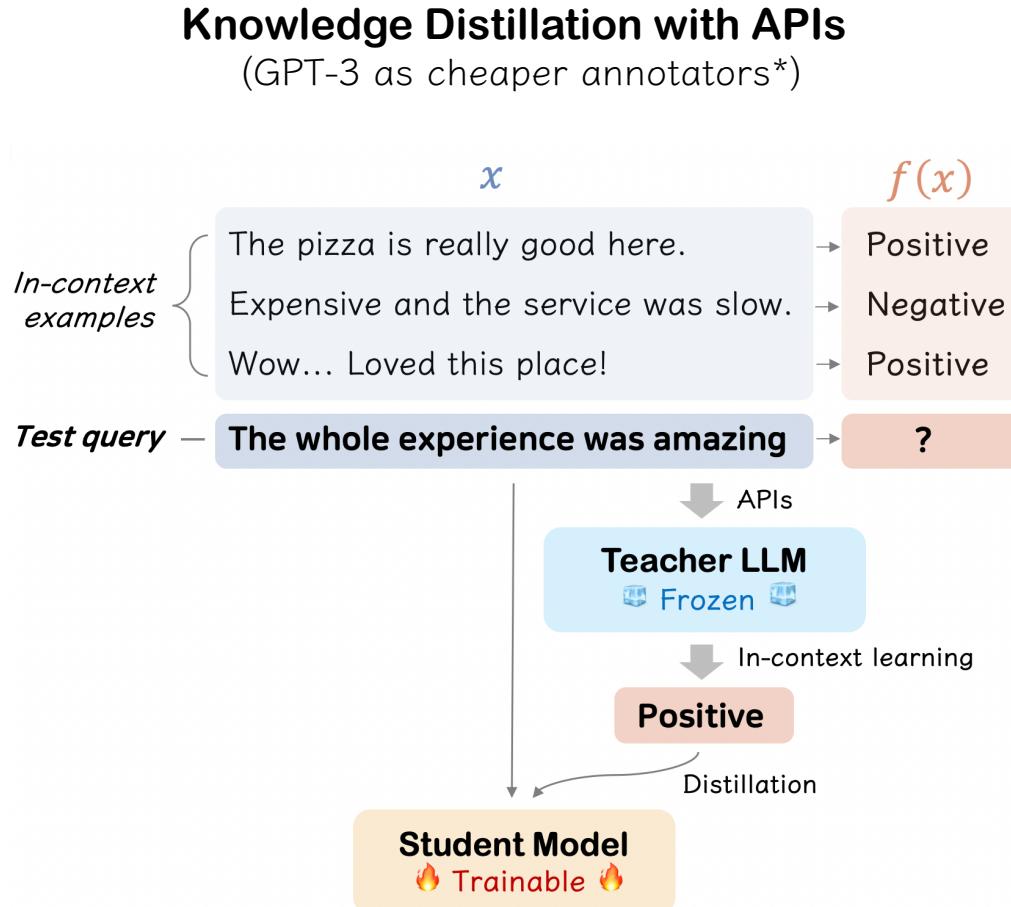
Knowledge Distillation with APIs

(GPT-3 as cheaper annotators*)

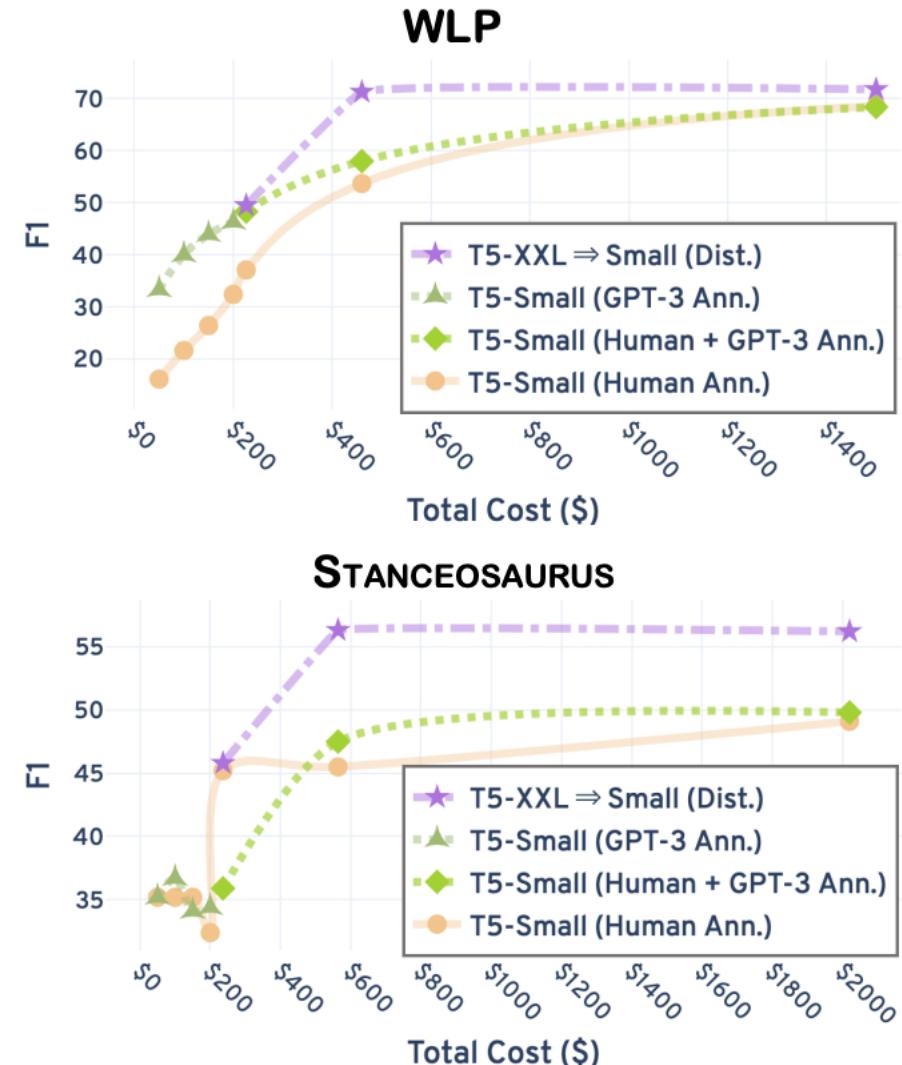


*Wang et al., 2022, Want To Reduce Labeling Cost? GPT-3 Can Help?

GPT-3.5 as an Annotator



*Wang et al., 2022, Want To Reduce Labeling Cost? GPT-3 Can Help?





Takeaways

*Q. Given a limited budget,
how to invest it to train a compact model
in an economically efficient manner?*

- ✓ In general, data annotation might not be the best practical solution in light of cost-efficiency; Scale up, then distill !



Takeaways

*Q. Given a limited budget,
how to invest it to train a compact model
in an economically efficient manner?*

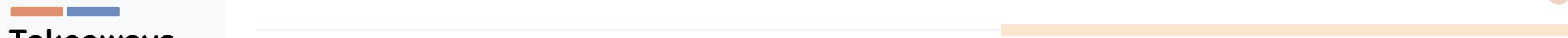
✓ In general, data annotation might not be the best practical solution in light of cost-efficiency; Scale up, then distill !

✓ For the best performance, however, data annotation is essential despite its inefficiency

Dist.: \$161 (81.0 F1) - max

Ann. : \$1,980 (81.0 F1)

\$17,443 (87.5 F1) - max



Takeaways

*Q. Given a limited budget,
how to invest it to train a compact model
in an economically efficient manner?*

✓ In general, data annotation might not be the best practical solution in light of cost-efficiency; Scale up, then distill !

✓ For the best performance, however, data annotation is essential despite its inefficiency

✓ Synthetic data generation using GPT-3.5 could be cost-efficient compared to humans, but still limited



Takeaways

*Q. Given a limited budget,
how to invest it to train a compact model
in an economically efficient manner?*

✓ In general, data annotation might not be the best practical solution in light of cost-efficiency; Scale up, then distill !

✓ For the best performance, however, data annotation is essential despite its inefficiency

✓ Synthetic data generation using GPT-3.5 could be cost-efficient compared to humans, but still limited

✓ More details and analyses in the paper, such as different sizes of large & compact models