Google DeepMind

# Many-Shot In-Context Learning

## NeurIPS, 2024 (Spotlight)

Rishabh Agarwal* , Avi Singh* , Lei M. Zhang† , Bernd Bohnet† , Luis Rosias† , Stephanie Chan† , Biao Zhang† , Ankesh Anand , Zaheer Abbas , Azade Nova , John D. Co-Reyes , Eric Chu , Feryal Behbahani , Aleksandra Faust and Hugo Larochelle
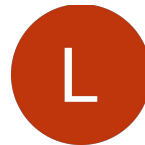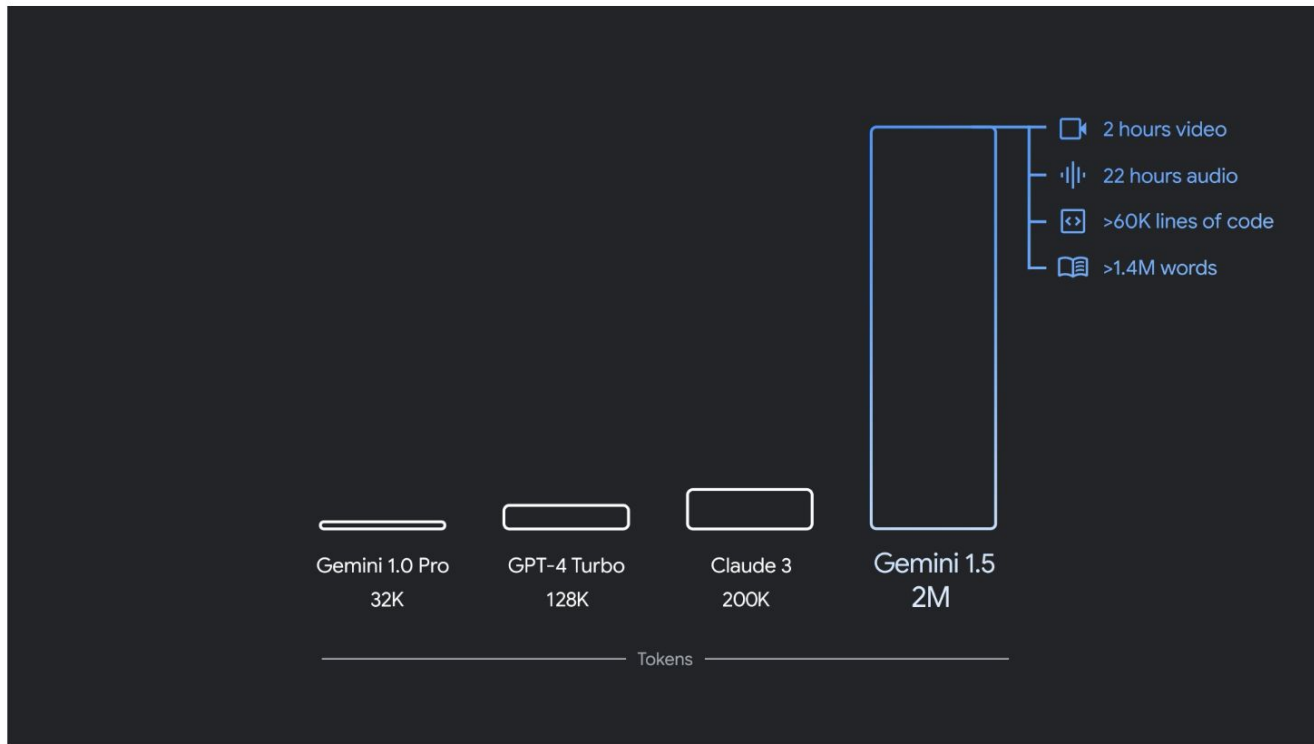
*Contributed equally, †Key contribution

# From few-shot to many-shot in-context learning (ICL)

1

# From few-shot to many-shot ICL



Context lengths of leading foundation models compared with Gemini 1.5's 2 million token capability
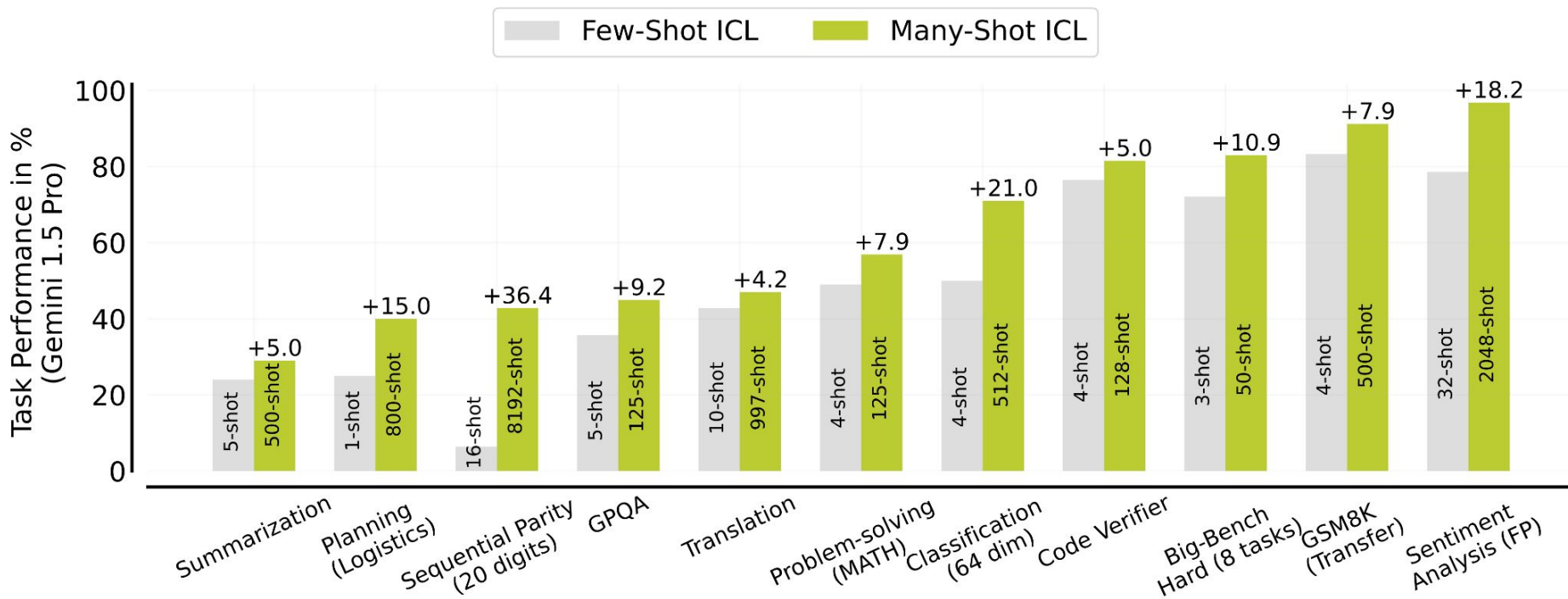
# From few-shot to many-shot ICL

## How many shots is "many-shot"?



Many-Shot ICL: Context Length versus Number of Shots

# From few-shot to many-shot ICL

## Does many-shot ICL improve performance? Yes!



Legend: Few-Shot ICL (grey), Many-Shot ICL (green)

Y-axis: Task Performance in % (Gemini 1.5 Pro), 0–100

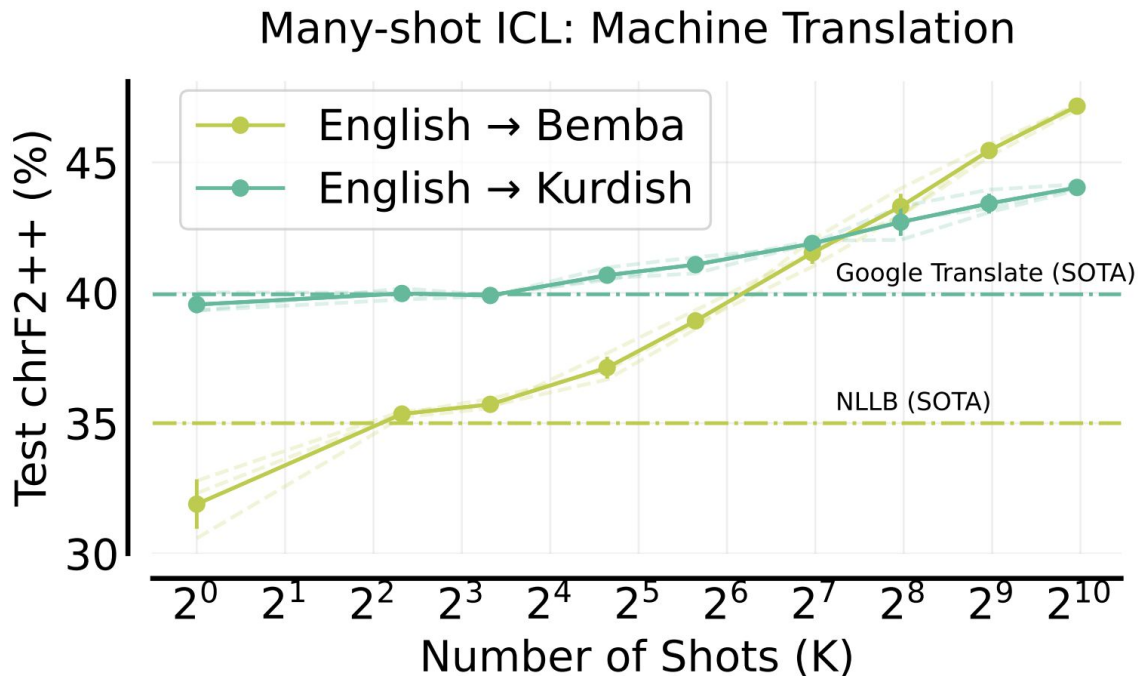| Task | Few-Shot | Many-Shot | Gain |
|------|----------|-----------|------|
| Summarization | 5-shot | 500-shot | +5.0 |
| Planning (Logistics) | 1-shot | 800-shot | +15.0 |
| Sequential Parity (20 digits) | 16-shot | 8192-shot | +36.4 |
| GPQA | 5-shot | 125-shot | +9.2 |
| Translation | 10-shot | 997-shot | +4.2 |
| Problem-solving (MATH) | 4-shot | 125-shot | +7.9 |
| Classification (64 dim) | 4-shot | 512-shot | +21.0 |
| Code Verifier | 4-shot | 128-shot | +5.0 |
| Big-Bench Hard (8 tasks) | 3-shot | 50-shot | +10.9 |
| GSM8K (Transfer) | 4-shot | 500-shot | +7.9 |
| Sentiment Analysis (FP) | 32-shot | 2048-shot | +18.2 |

# Many-shot ICL examples

2

# Machine translation on low-resource languages

## Beating SOTA systems using many-shot ICL.
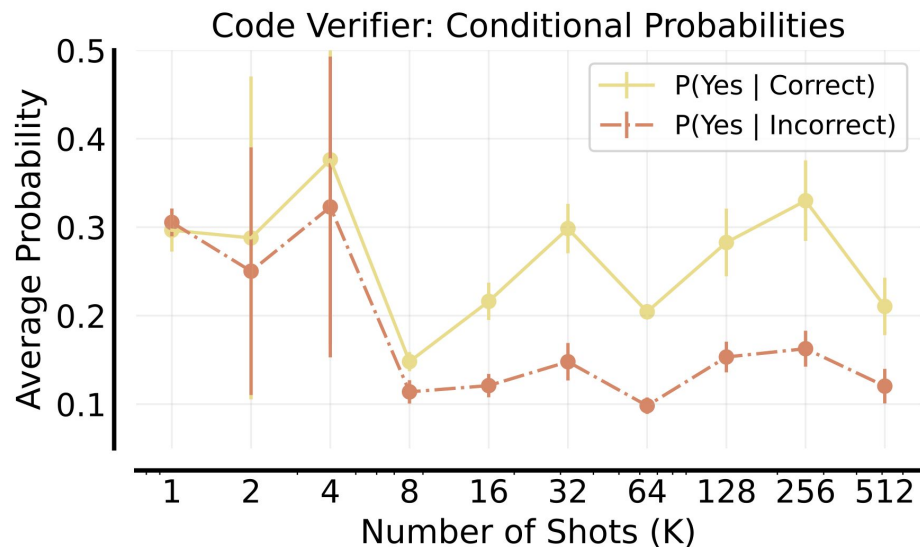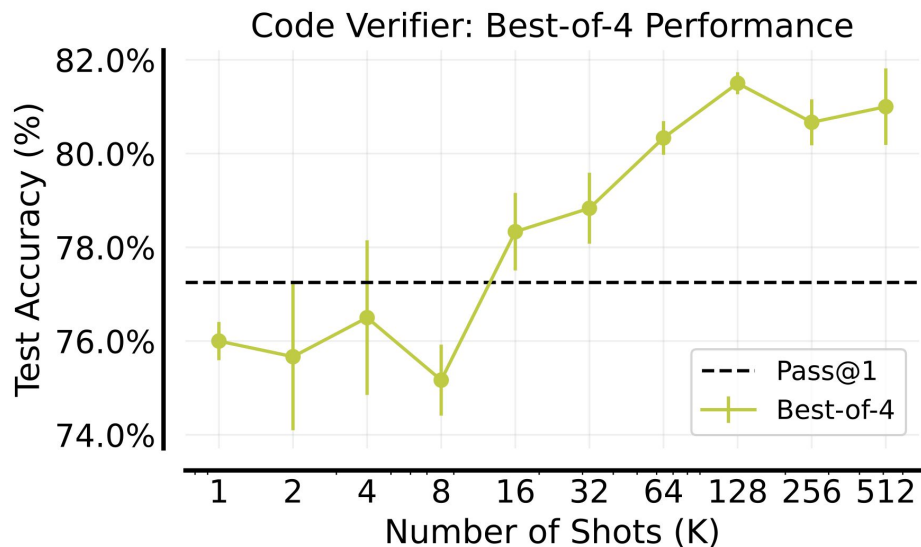


Many-shot ICL: Machine Translation

# Logistics Planning

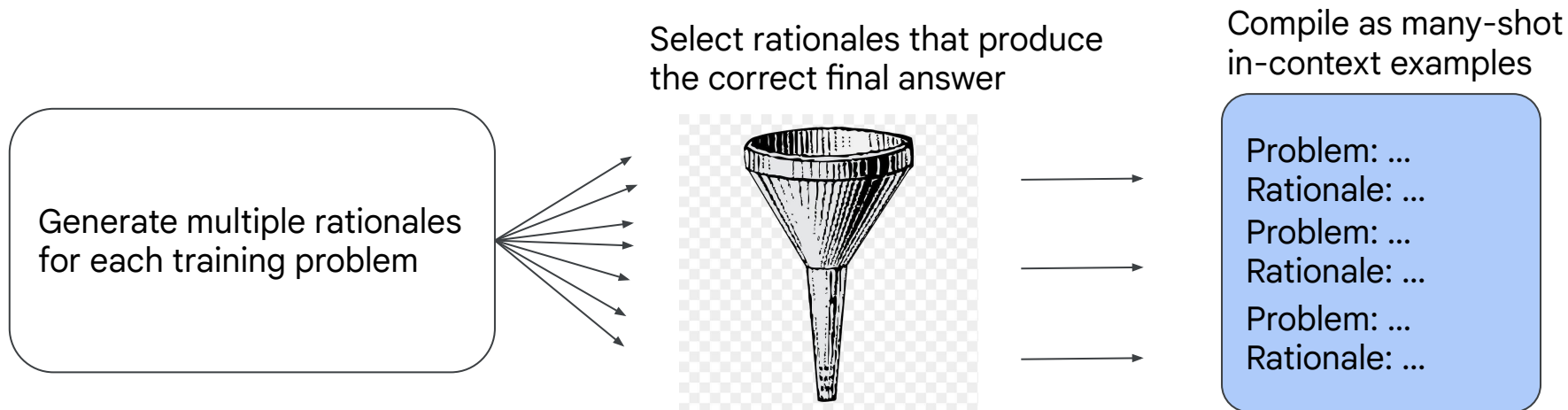# Code verifier

**Code verifier without fine-tuning!**

**Many more examples of effective many-shot ICL: e.g., Planning, Summarization. See paper for details!**

# Many-shot learning without human-written rationales

3

# Human-written rationales or demonstrations can be expensive to collect... can we do without?

## Reinforced ICL: use model-generated rationales

Generate multiple rationales for each training problem

Select rationales that produce the correct final answer



Compile as many-shot in-context examples

Problem: ...
Rationale: ...
Problem: ...
Rationale: ...
Problem: ...
Rationale: ...

*called "reinforced" because of equivalence to expectation-maximization RL algorithm

# Human-written rationales or demonstrations can be expensive to collect... can we do without?

## Unsupervised ICL: get rid of rationales/solutions entirely!

**Preamble**

You will be provided Problems similar to the ones below:

**Long list of unsolved problems**

Problem: ...
Problem: ...
Problem: ...

*Many-shot to teach the problem space*

**Instruction**

Now, I am going to give you a series of demonstrations of math Problems and Solutions. When you respond, respond only with the Solution of the final Problem, thinking step by step.
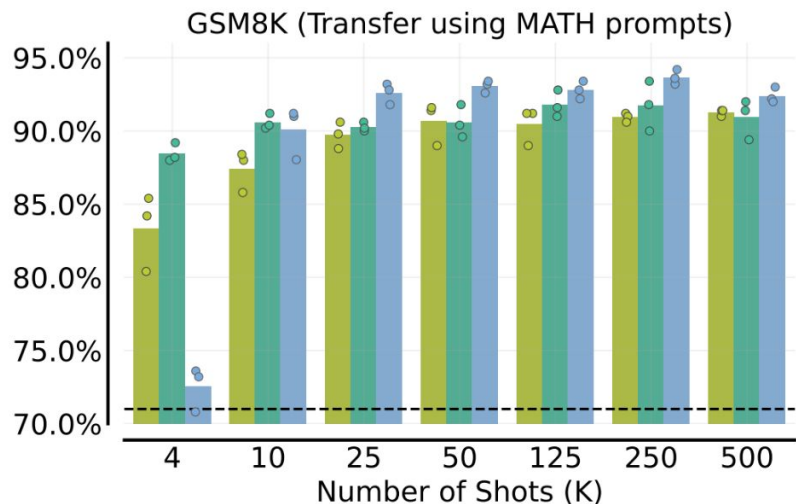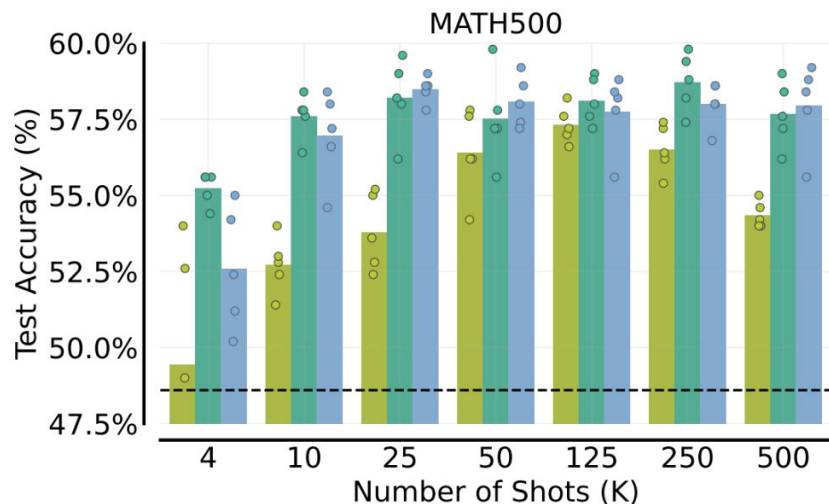
**Short list of problems with solutions**

Problem: ...
Solution: ...
Problem: ...

*Few-shot to teach the format*

# Problem-Solving: Hendrycks MATH & GSM8K

**Reinforced and Unsupervised ICL can outperform ICL with human-written solutions!**

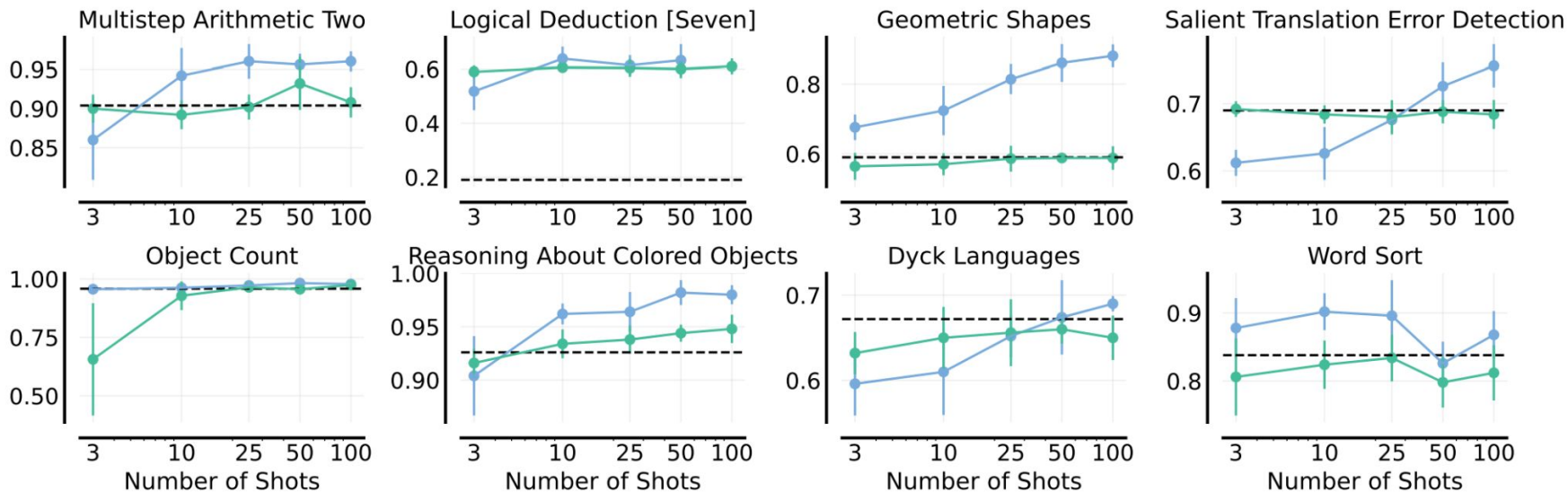The Hendrycks MATH prompts transfer well to GSM8k (another math dataset)



MATH500

GSM8K (Transfer using MATH prompts)

Legend:
- `-----` 4-shot InnerMono. MATH Prompt
- ICL (Ground-Truth) (human-written)
- Unsupervised ICL
- Reinforced ICL

# Algorithmic and Symbolic Reasoning: Big-Bench Hard

**Generally: Reinforced > Unsupervised > Human-written**
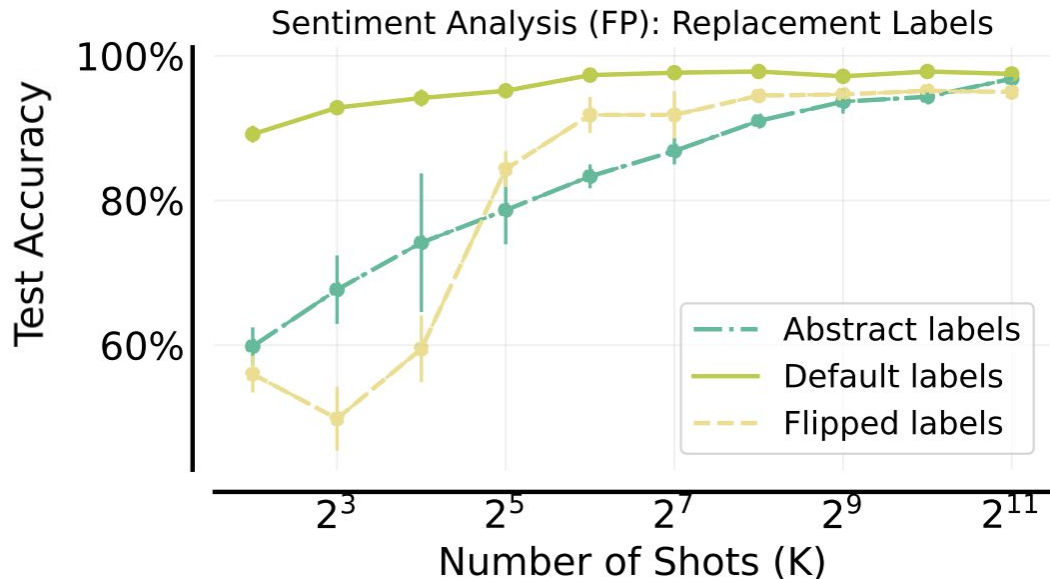
● with greater improvement for more shots2

# Analyzing many-shot ICL

4

# Many-shot ICL can overcome pre-training biases

Previous work (Kossen et al, 2023) suggest that ICL has difficulty unlearning biases derived from pre-training data…



Sentiment Analysis (FP): Replacement Labels

| Default (original) | Flipped (rotated) | Abstract |
|---|---|---|
| negative | neutral | A |
| neutral | positive | B |
| positive | negative | C |

…but with enough shots, new labels eventually approach performance of original labels

# High-dimensional functions:
# Binary Linear Classification in High Dimensions

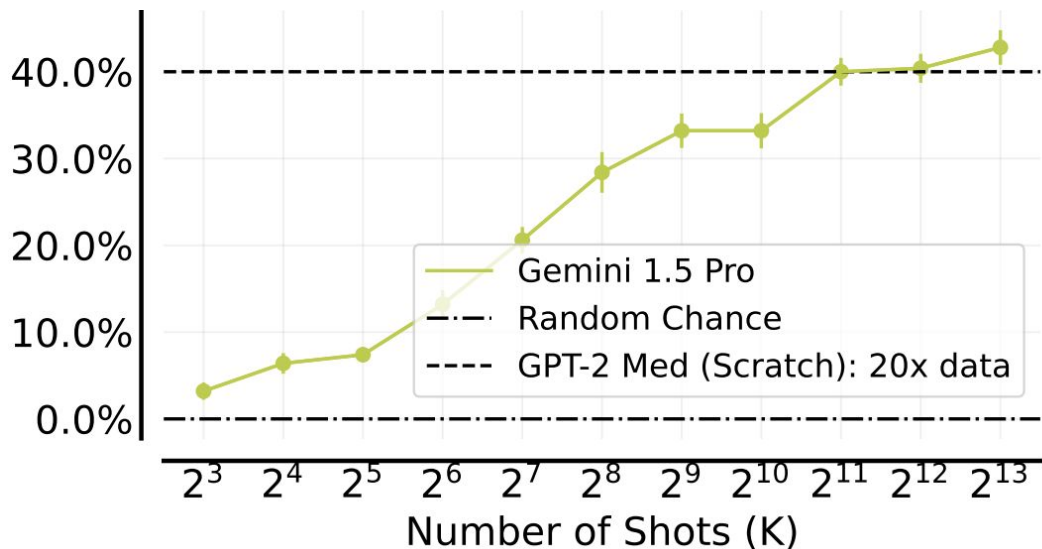Many-shot ICL nearly matches strong baseline (k-nearest neighbors)

# High-dimensional functions:
# Sequential Parity (20 digits)

Does the binary input sequence so far contain even or odd number of 1s?

**Input:** 1 0 1 1 0 0 0 1 1 1 0 0 0 0 1 0 0 1 1 1
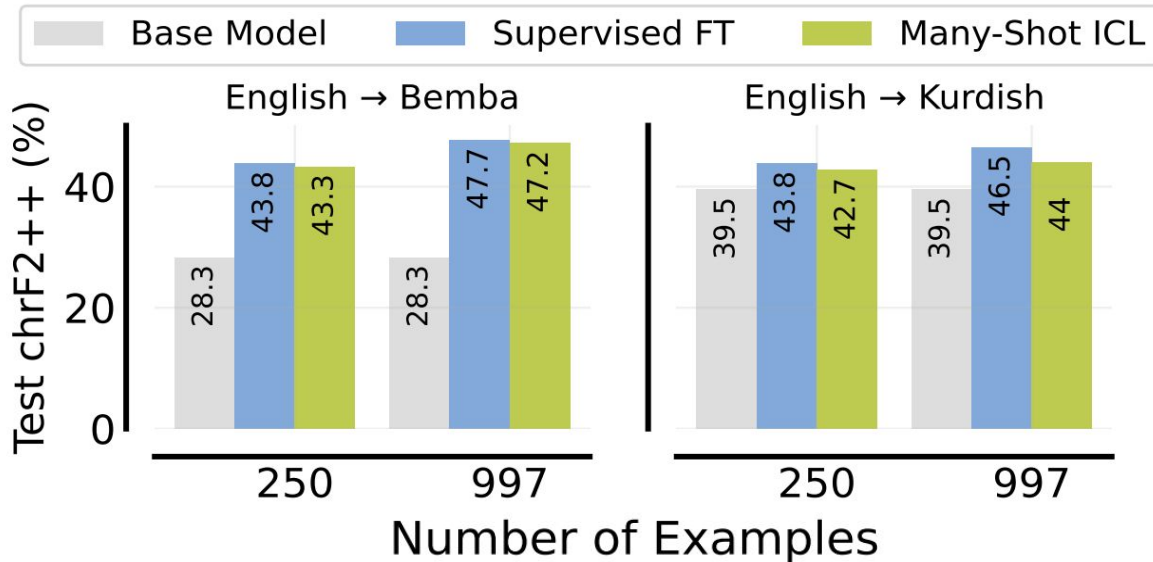**Label:** Odd Odd Even Odd Odd Odd Odd Even Odd Even Even Even Even Even Odd Odd Odd Even Odd Even

believed to be a fundamental limitation of self-attention (Chiang and Cholak, 2022)
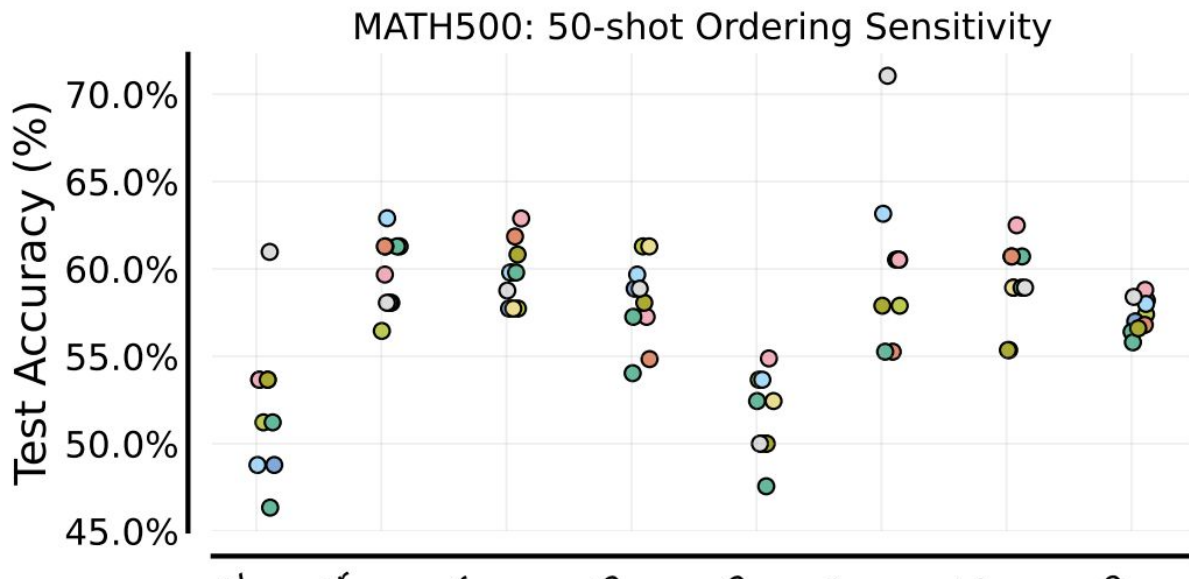
Many-shot ICL improves monotonically until 2^13 examples!



And outperforms a GPT-2 sized model trained from scratch on 20x more data

# Many-shot ICL can have similar performance to SFT (translation task)



- ICL has no training cost but potentially higher inference cost (can mitigate with context caching)
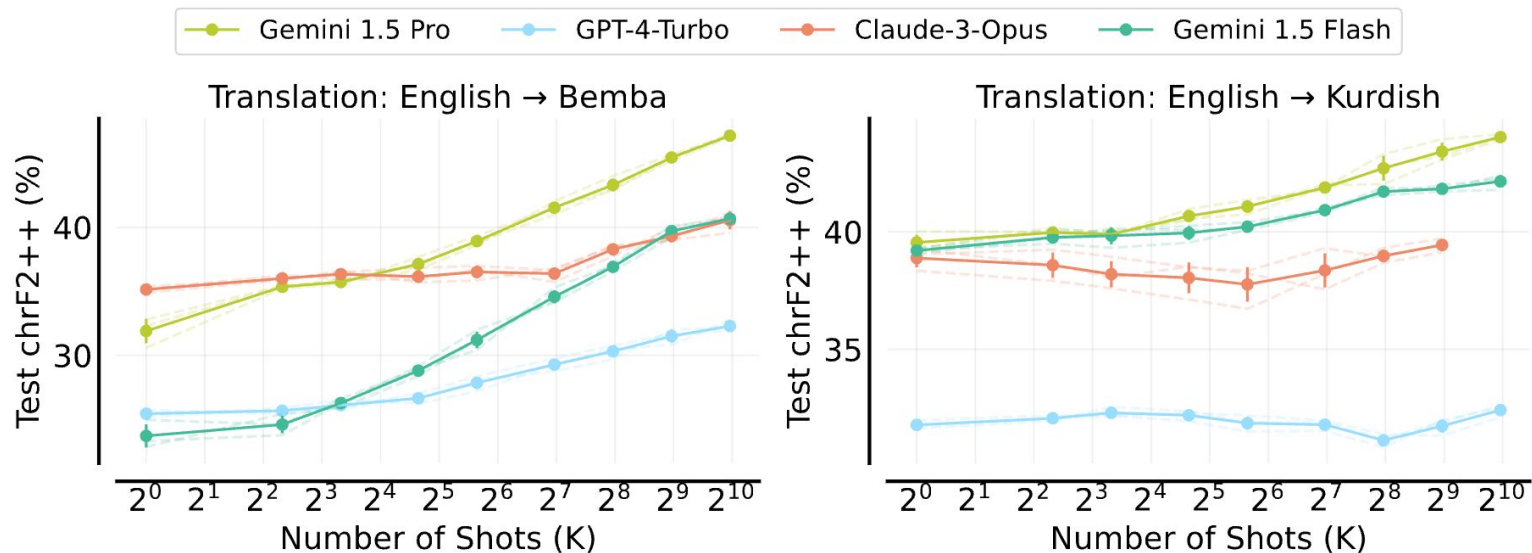
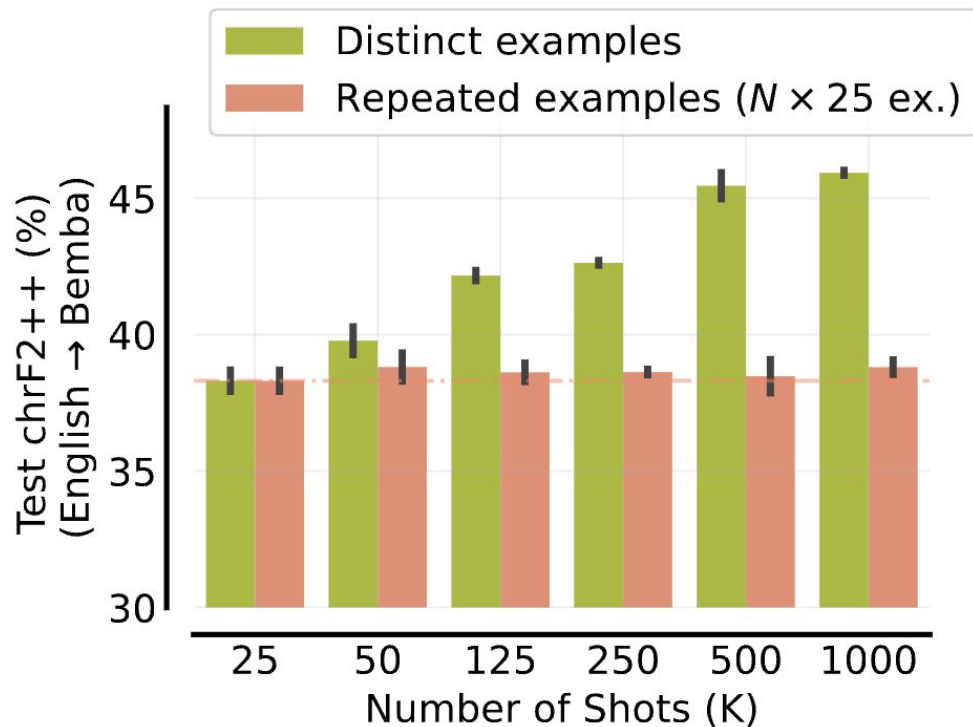# Many-shot ICL can be sensitive to example ordering



MATH500: 50-shot Ordering Sensitivity

Each colored point is a different ordering

# Comparison of frontier models



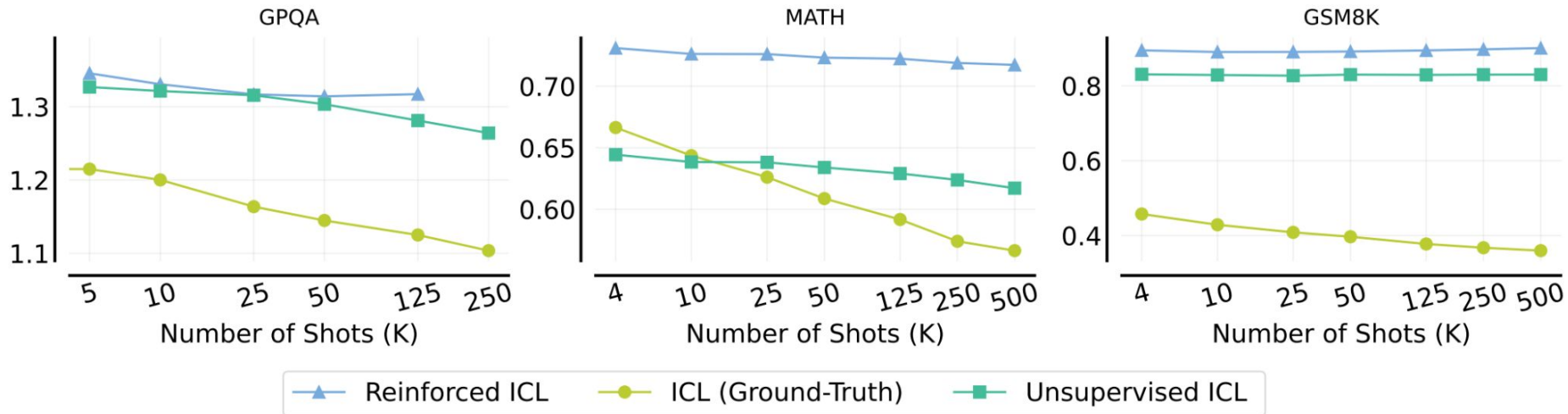Varying levels of benefit from many-shot ICL

# Increasing Context Length or More Information?



Many-shot performance with distinct examples vs repeating the same 25 examples N times on low-resource MT.

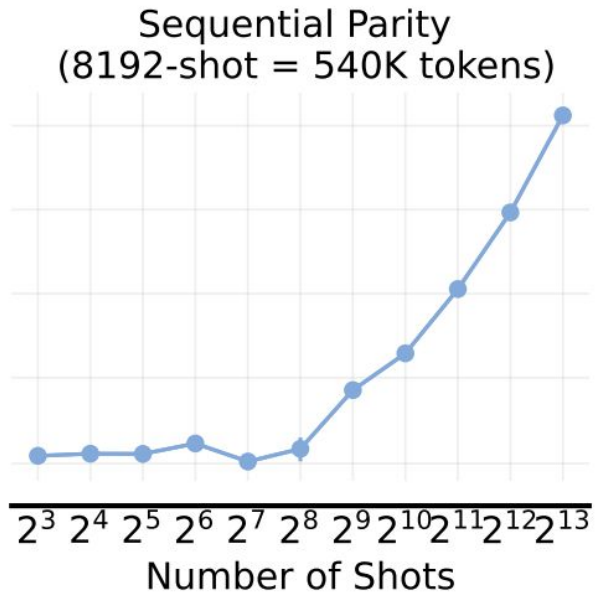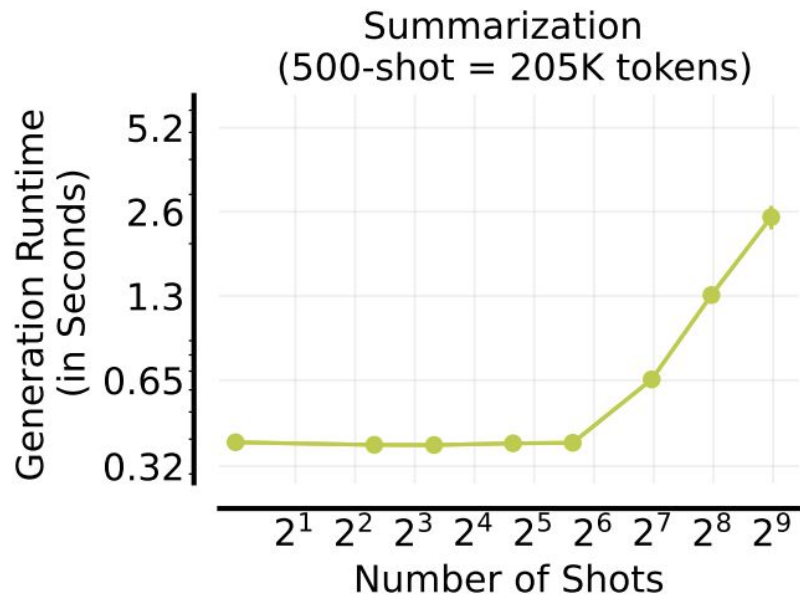# Long-context scaling laws may not predict ICL performance



Negative Log-Likelihood on Ground-Truth Solutions

GPQA  MATH  GSM8K

Number of Shots (K)

Reinforced ICL — ICL (Ground-Truth) — Unsupervised ICL

NLL is not a reliable predictor of ICL performance

- NLL consistently decreases, even though ICL worsens beyond 125 shots
- NLL for human-written rationales is lower than for model-written rationales, even though actual performance is often worse

# Inference costs



**Summarization**
(500-shot = 205K tokens)

**Sequential Parity**
(8192-shot = 540K tokens)

**Summary**

1. **Many-shots can improve performance up to 1000s of shots**
   Long-context models enable this (to varying degrees)

2. **Model-generated or unsupervised prompts** can often outperform human-written prompts

3. **Analyses**:
   - Many-shot ICL can overcome pretraining biases
   - Many-shot ICL can have similar performance to SFT
   - NLL is not a reliable predictor of ICL performance

**Q:** What are the mechanisms underlying many-shot learning? Why do particular tasks benefit more?

**Q:** Why does performance sometimes degrade after many shots?

**Q:** Why does example ordering matter?