## Are Emergent Abilities of Large Language Models a Mirage?

Rylan Schaeffer, Brando Miranda, and Sanmi Koyejo

Computer Science, Stanford University

### Abstract

Recent work claims that large language models display *emergent abilities*, abilities not present in smaller-scale models that are present in larger-scale models. What makes emergent abilities intriguing is two-fold: their *sharpness*, transitioning seemingly instantaneously from not present to present, and their *unpredictability*, appearing at seemingly unforeseeable model scales. Here, we present an alternative explanation for emergent abilities: that for a particular task and model family, when analyzing fixed model outputs, emergent abilities appear due the researcher's choice of metric rather than due to fundamental changes in model behavior with scale. Specifically, nonlinear or discontinuous metrics produce apparent emergent abilities, whereas linear or continuous metrics produce smooth, continuous, predictable changes in model performance. We present our alternative explanation in a simple mathematical model, then test it in three complementary ways: we (1) make, test and confirm three predictions on the effect of metric choice using the InstructGPT/GPT-3 family on tasks with claimed emergent abilities, (2) make, test and confirm two predictions about metric choices in a meta-analysis of emergent abilities on BIG-Bench; and (3) show how to choose metrics to produce never-before-seen seemingly emergent abilities in multiple vision tasks across diverse deep networks. Via all three analyses, we provide evidence that alleged emergent abilities evaporate with different metrics or with better statistics, and may not be a fundamental property of scaling AI models.

# Are Emergent Abilities of LLMs a Mirage?

By Rylan Schaeffer, Brando Miranda, and Sanmi Koyejo

Presented by Chao Chen (Michelle)

1

# Emergent Abilities

**Sparks of Artificial General Intelligence:**
**Early experiments with GPT-4**

Sébastien Bubeck    Varun Chandrasekaran    Ronen Eldan    Johannes Gehrke
Eric Horvitz    Ece Kamar    Peter Lee    Yin Tat Lee    Yuanzhi Li    Scott Lundberg
Harsha Nori    Hamid Palangi    Marco Tulio Ribeiro    Yi Zhang

Microsoft Research

**Abstract**

been developing and refining large language models (LLMs)
variety of domains and tasks, challenging our understanding
eveloped by OpenAI, GPT-4 [Ope23], was trained using an
this paper, we report on our investigation of an early version
nt by OpenAI. We contend that (this early version of) GPT-
th ChatGPT and Google's PaLM for example) that exhibit
odels. We discuss the rising capabilities and implications of
ts mastery of language, GPT-4 can solve novel and difficult
medicine, law, psychology and more, without needing any
asks, GPT-4's performance is strikingly close to human-level
models such as ChatGPT. Given the breadth and depth of
easonably be viewed as an early (yet still incomplete) version

**Predictability and Surprise in Large Generative Models**

DEEP GANGULI*, DANNY HERNANDEZ*, LIANE LOVITT*, NOVA DASSARMA†, T
ANDY JONES†, NICHOLAS JOSEPH†, JACKSON KERNION†, BEN MANN†, AMA
YUNTAO B
SHOWK, S
NEEL NAN
JARED KAI
Anthropic, US

Large-scale pre
Megatron-Turi
the policy impli
training distribu
the high-level p
qualities make
can lead to soci
experiments to
combine to give
conclude with a
impact. We inte
about the poten
want to analyze

ACM Referen
Deep Ganguli,
Ben Mann, Am

# Language Models are Few-Shot Learners

Tom B. Brown*        Benjamin Mann*        Nick Ryder*        Melanie Subbiah*

Jared Kaplan†    Prafulla Dhariwal    Arvin

Amanda Askell    Sandhini Agarwal    Ariel H

Rewon Child    Aditya Ramesh    Daniel

Christopher Hesse    Mark Chen    E

Benjamin Chess        Jack

# Emergent Abilities of Large Language Models

Jason Wei [1]                                          jasonwei@google.com
Yi Tay [1]                                              yitay@google.com
Rishi Bommasani [2]                                     nlprishi@stanford.edu
Colin Raffel [3]                                        craffel@gmail.com
Barret Zoph [1]                                         barretzoph@google.com

# Emergent Abilities are a Mirage

*Sharp and unpredictable changes are induced by researcher's choice of metric. Model family's per-token error rate changes smoothly, continuously, and predictable*

Part I:    Intuition on emergent abilities

Part II:    Empirically prove hypothesis InstructGPT/GPT-3
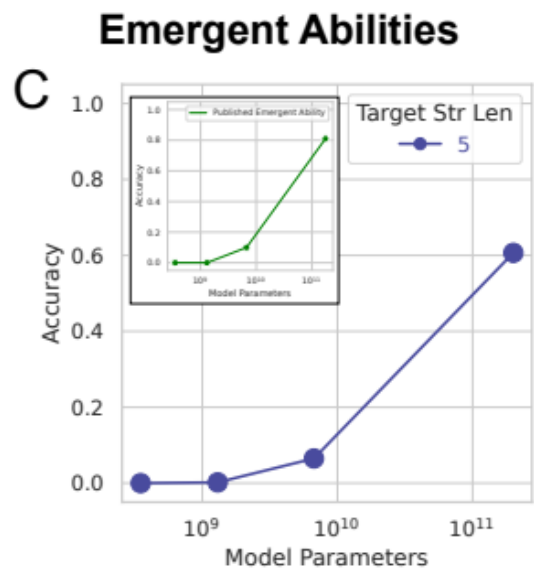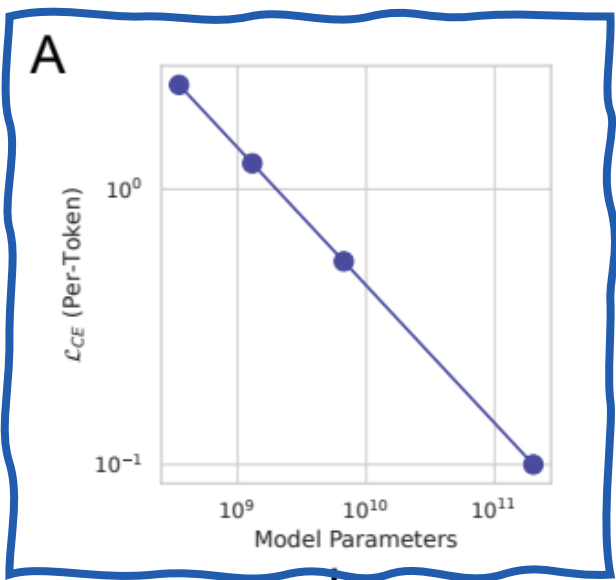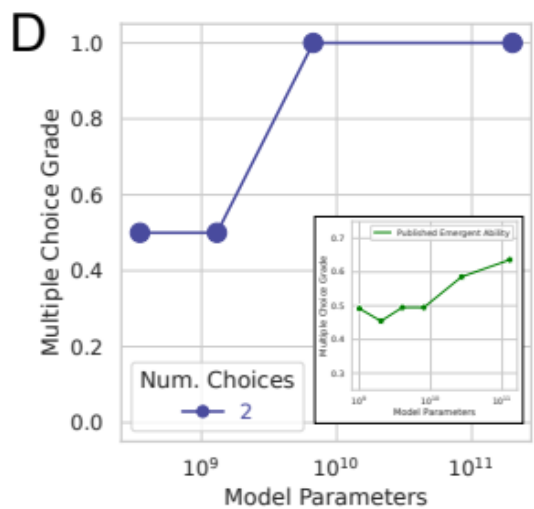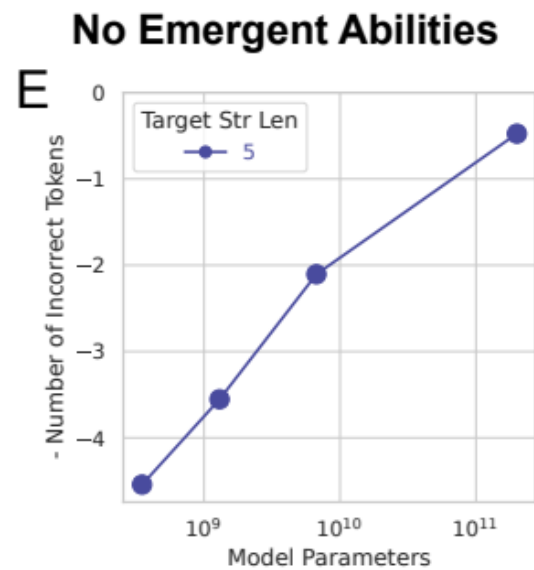
Part III:    Meta-analysis of emergent abilities

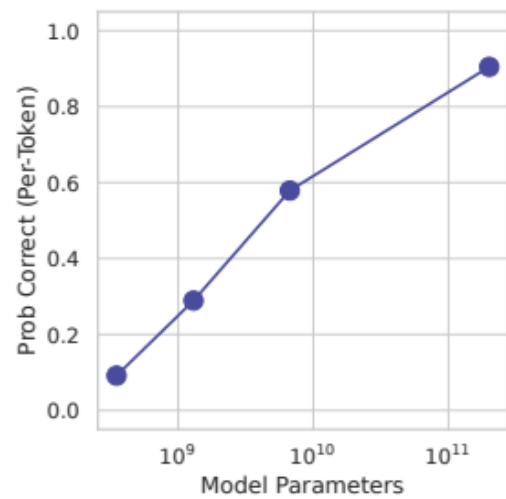Part IV:    Induce emergent abilities on vision models

Part I:     Intuition on emergent abilities

Part II:    Empirically prove hypothesis InstructGPT/GPT-3

Part III:   Meta-analysis of emergent abilities

Part IV:   Induce emergent abilities on vision models

**Emergent Abilities**

C

Target Str Len — 5

Published Emergent Ability

**Nonlinearly score LLM outputs**

A

$\mathcal{L}_{CE}$ (Per-Token)

Model Parameters

**No Emergent Abilities**

E

Target Str Len — 5

**Linearly score LLM outputs**

B

$$p(\text{single token correct}) = \exp\left(-\mathcal{L}_{CE}(N)\right)$$

Prob Correct (Per-Token)

Model Parameters

D

Num. Choices — 2

Published Emergent Ability

**Discontinuously score LLM outputs**

**Continuously score LLM outputs**

F

Num. Choices — 2

A

$$\mathcal{L}_{CE}(N) = \left(\frac{N}{C}\right)^{\alpha}$$

Cross Entropy Loss

$$\mathcal{L}_{CE}(N) := -\sum_{v \in V} p(v) \log \hat{p}_N(v)$$

One-hot distribution

$$\mathcal{L}_{CE}(N) = -\log \hat{p}_N(v^*)$$

**Emergent Abilities**

C

Nonlinearly score LLM outputs

A

B

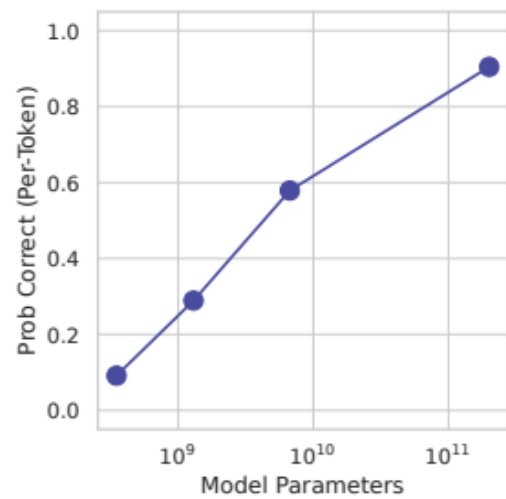$$p(\text{single token correct}) = \exp\left(-\mathcal{L}_{CE}(N)\right)$$

D

Discontinuously score LLM outputs

**No Emergent Abilities**

E

Linearly score LLM outputs

Continuously score LLM outputs

F

8

$$p(\text{single token correct}) = \exp\left(-\boxed{\mathcal{L}_{CE}(N)}\right)$$



$$\mathcal{L}_{CE}(N) = -\log \hat{p}_N(v^*)$$

$$p_N(\text{single token correct}) = \exp\left(-\boxed{\left(\frac{N}{C}\right)^{\alpha}}\right)$$

**Emergent Abilities**

**No Emergent Abilities**

A

$\mathcal{L}_{CE}$ (Per-Token)

Model Parameters

C

Accuracy

Published Emergent Ability

Target Str Len
5

Accuracy

Model Parameters

Model Parameters

**Nonlinearly score LLM outputs**

E

- Number of Incorrect Tokens

Target Str Len
5

Model Parameters

**Linearly score LLM outputs**

B

$$p(\text{single token correct}) = \exp\left(-\mathcal{L}_{CE}(N)\right)$$

Prob Correct (Per-Token)

Model Parameters

D

Multiple Choice Grade

Num. Choices
2

Published Emergent Ability

Multiple Choice Grade

Model Parameters

Model Parameters

**Discontinuously score LLM outputs**

**Continuously score LLM outputs**

F

- Brier Score

Num. Choices
2

Model Parameters

10

# Emergent Abilities
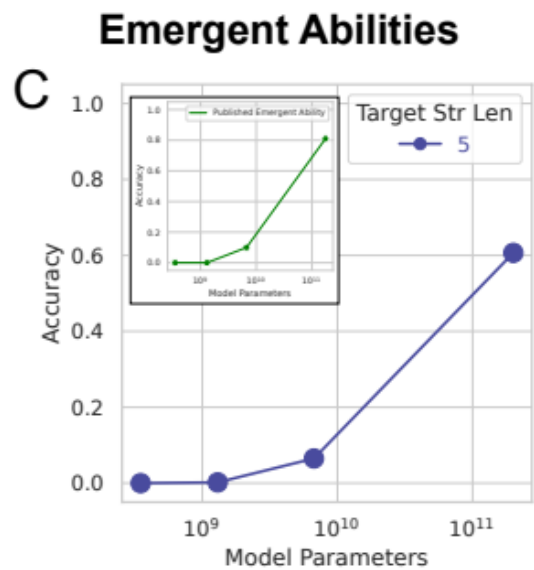


$$\text{Accuracy}(N)$$
$$\approx p_N(\text{single token correct})^{\text{num of tokens}}$$
$$= \exp(-\left(\frac{N}{C}\right)^{\alpha})^{L}$$

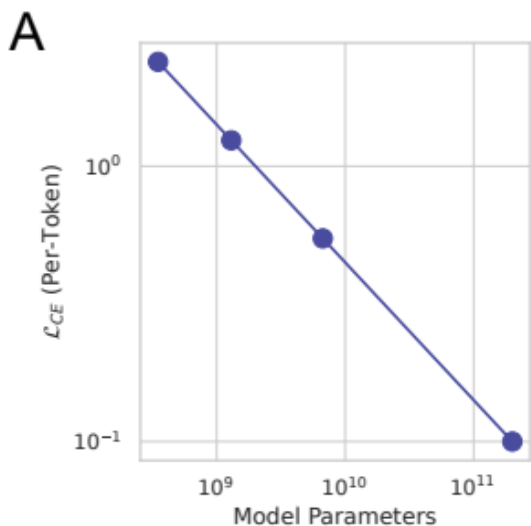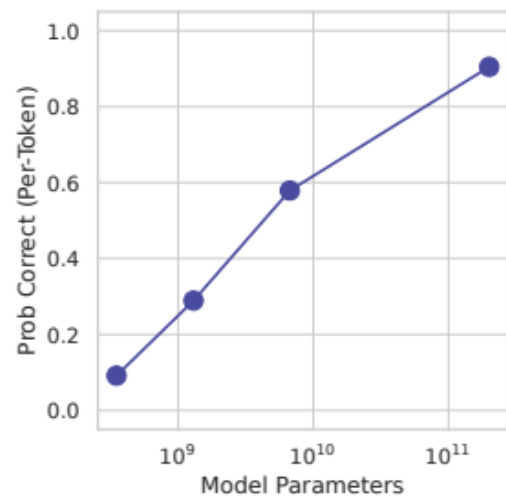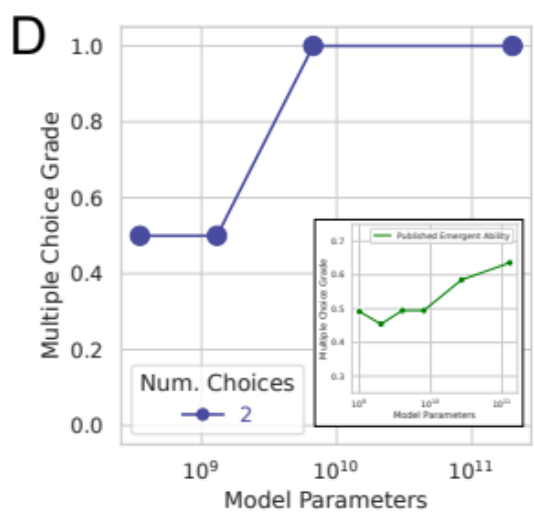# No Emergent Abilities



$$\text{Token Edit Distance }(N)$$
$$\approx L \cdot \left(1 - p_N(\text{single token correct})\right)$$
$$= L \cdot (1 - \exp(-\left(\frac{N}{C}\right)^{\alpha}))$$

**Emergent Abilities**

C — Accuracy vs Model Parameters, Target Str Len 5

**No Emergent Abilities**

E — - Number of Incorrect Tokens vs Model Parameters, Target Str Len 5

A — $\mathcal{L}_{CE}$ (Per-Token) vs Model Parameters

**Nonlinearly score LLM outputs**

**Linearly score LLM outputs**

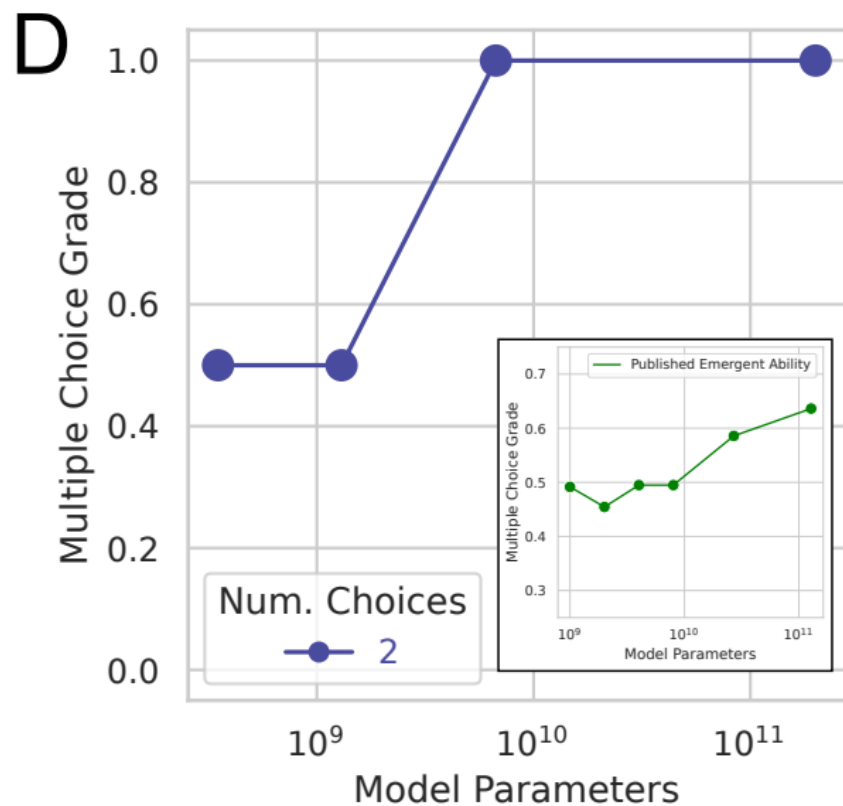B — $p(\text{single token correct}) = \exp\left(-\mathcal{L}_{CE}(N)\right)$ — Prob Correct (Per-Token) vs Model Parameters

D — Multiple Choice Grade vs Model Parameters, Num. Choices 2

**Discontinuously score LLM outputs**

**Continuously score LLM outputs**

F — - Brier Score vs Model Parameters, Num. Choices 2

$$\text{Multiple Choice Grade}(N)$$
$$\approx \sum_{i=0}^{n} 1_{[p(v*) > p(v)]}$$

$$\text{Brier Score}(N)$$
$$\approx \frac{1}{n} \sum_{i=0}^{n} (\hat{p}(v^*) - 1_{[v*]})^2$$

Part I:    Intuition on emergent abilities

Part II:    Empirically prove hypothesis InstructGPT/GPT-3

Part III:   Meta-analysis of emergent abilities

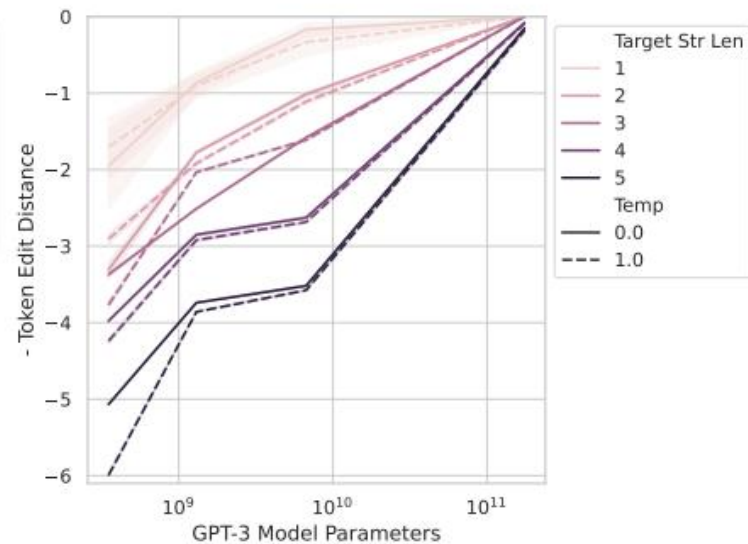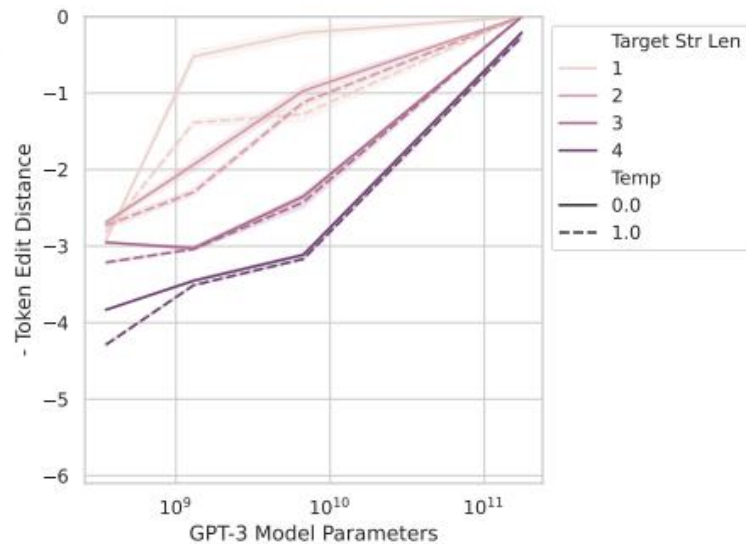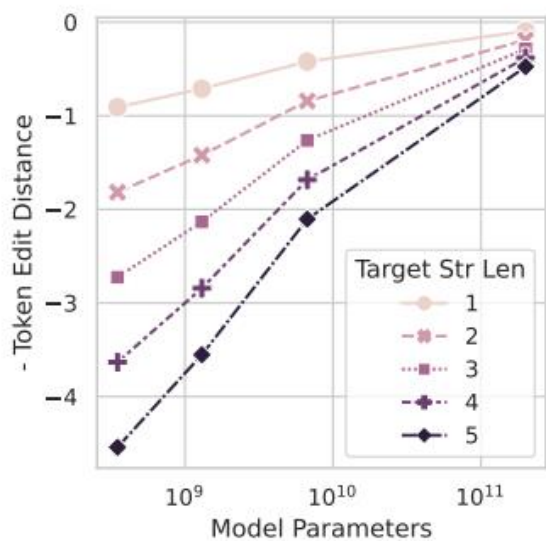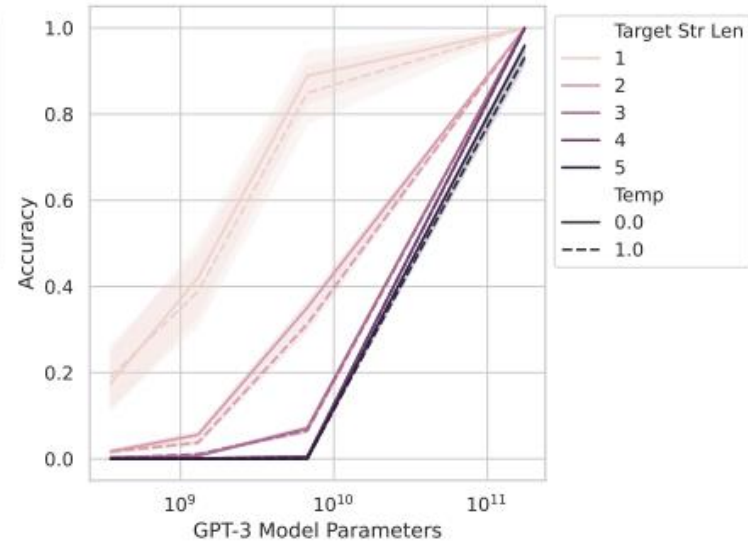Part IV:  Induce emergent abilities on vision models

Part I:      Intuition on emergent abilities

Part II:    Empirically prove hypothesis InstructGPT/GPT-3

Part III:   Meta-analysis of emergent abilities
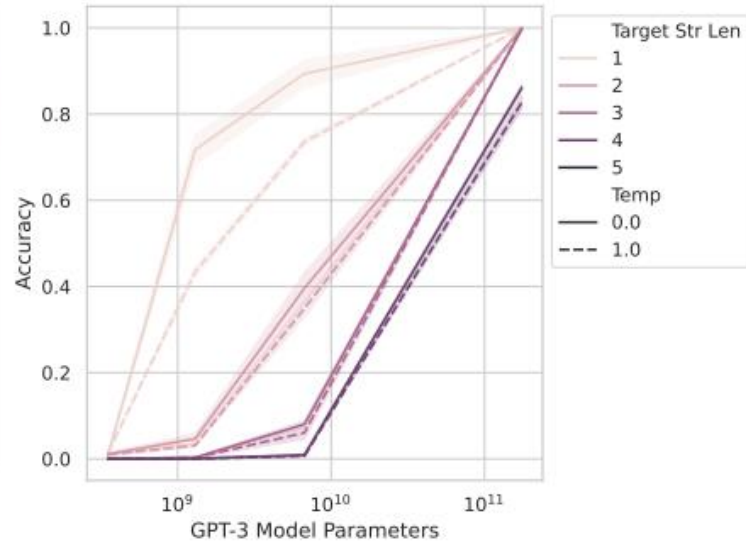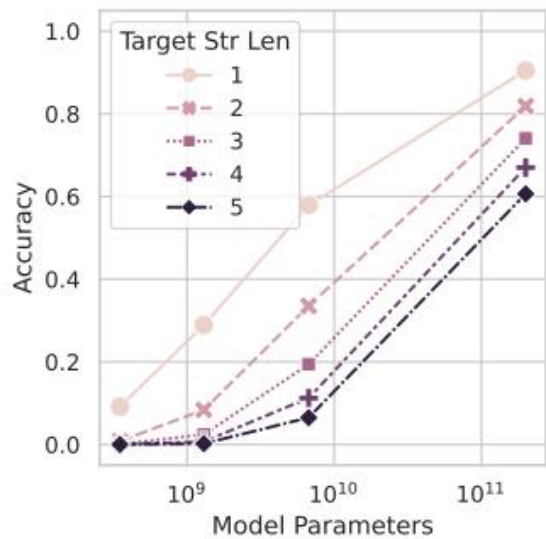
Part IV:  Induce emergent abilities on vision models

# Predictions

1. Emergent abilities disappear with different metrics.
2. Emergent abilities disappear with better statistics.

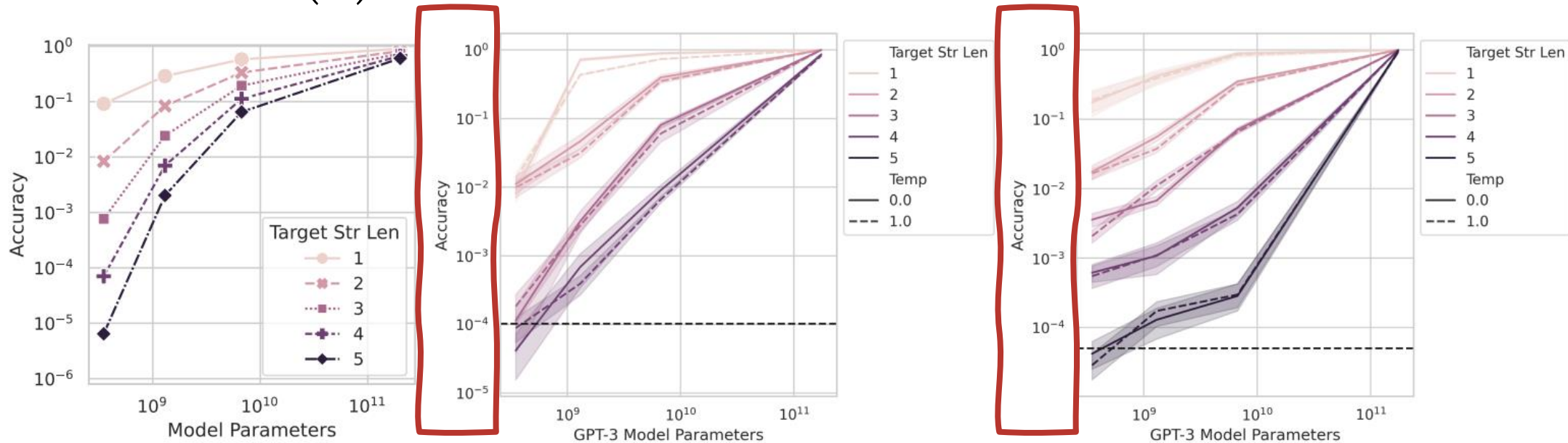$$\mathcal{L}_{CE}(N) = \left(\frac{N}{C}\right)^{\alpha}$$

Two-digit multiplication

Four-digit addition

$$\mathcal{L}_{CE}(N) = \left(\frac{N}{C}\right)^{\alpha}$$

# Two-digit multiplication

# Four-digit addition

Part I:    Intuition on emergent abilities

Part II:   Empirically prove hypothesis InstructGPT/GPT-3

Part III:  Meta-analysis of emergent abilities

Part IV:  Induce emergent abilities on vision models

Part I:    Intuition on emergent abilities

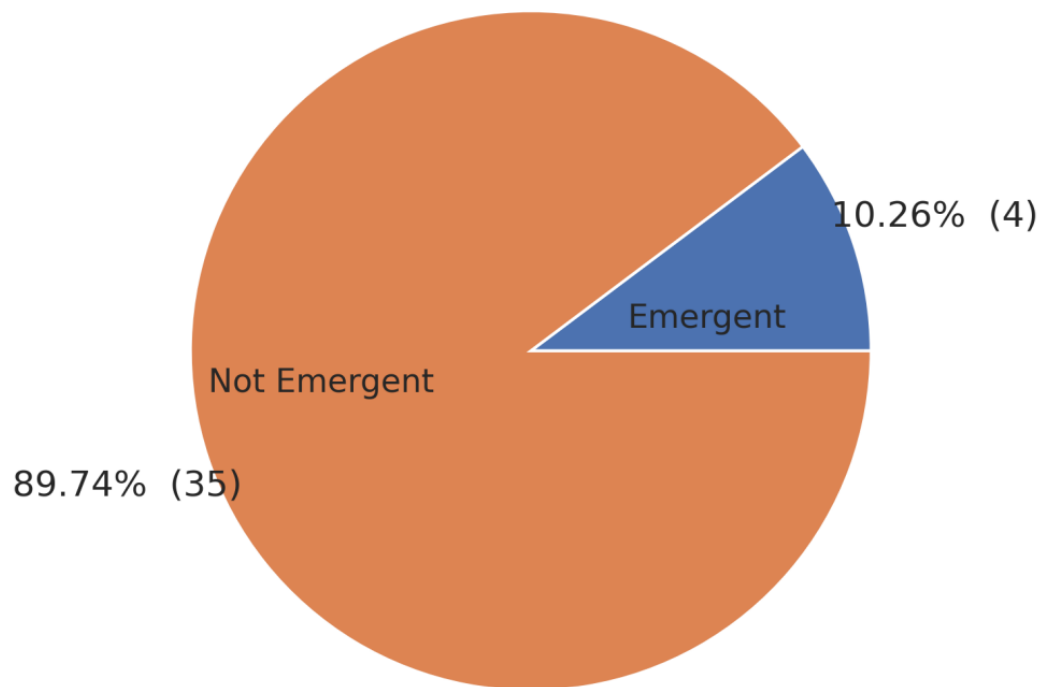Part II:   Empirically prove hypothesis InstructGPT/GPT-3

Part III:  Meta-analysis of emergent abilities

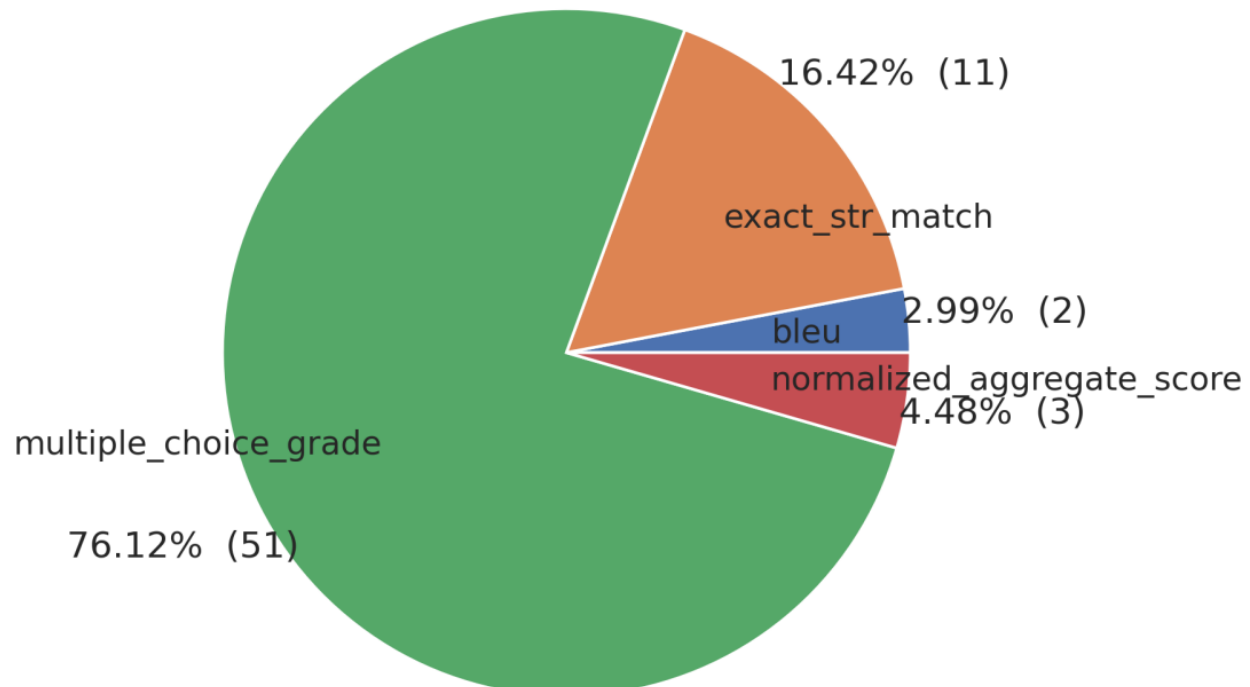Part IV:  Induce emergent abilities on vision models
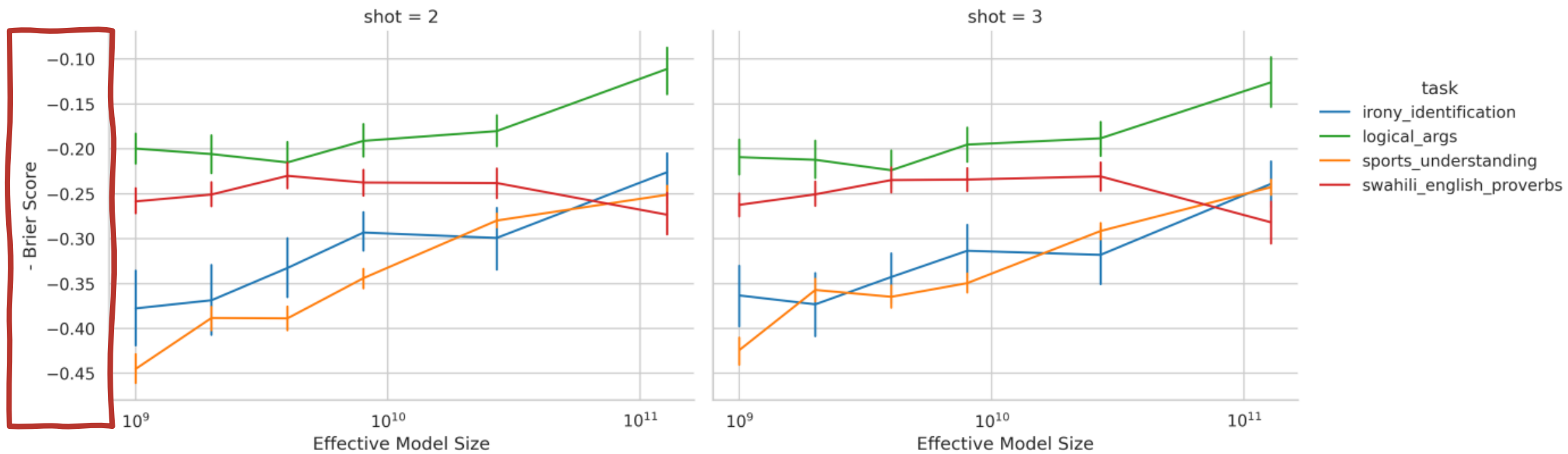
# Predictions

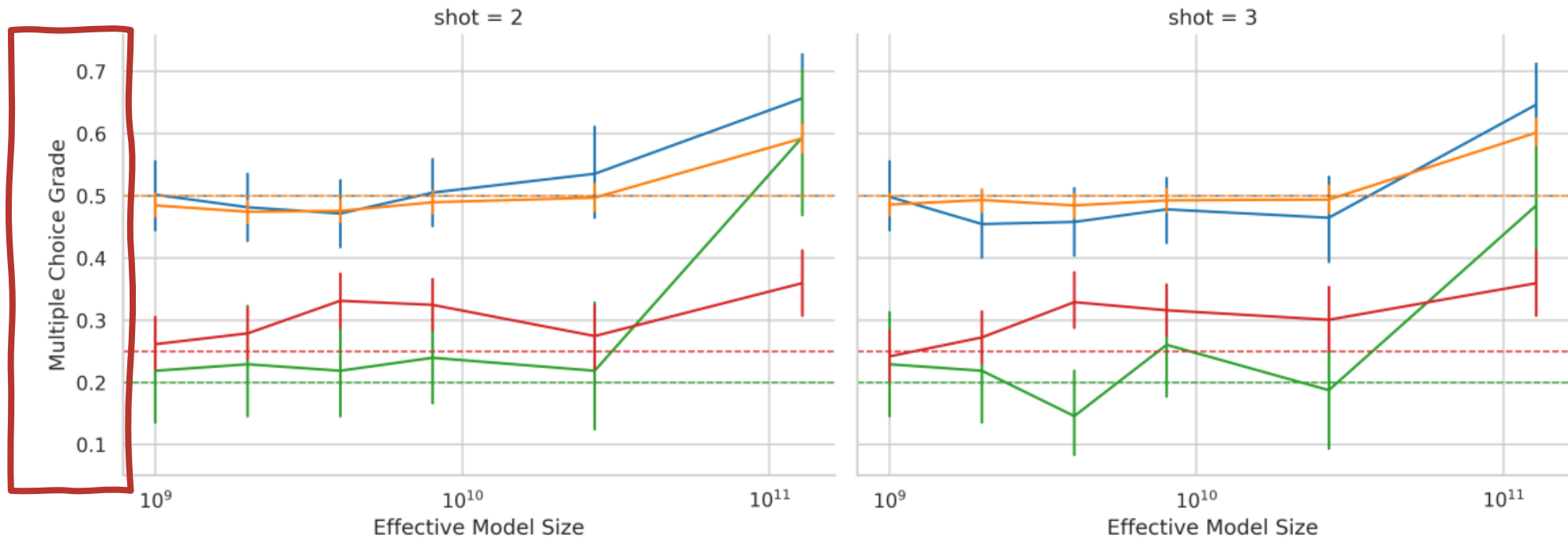- Emergent abilities appear with discontinuous/non-linear metrics.
- Emergent abilities disappear after changing metric.

% of Metrics with > 1 Model-Task
Pair Exhibiting Emergent Abilities

Metrics of Model-Task Pairs
Exhibiting Emergent Abilities

Top row (shot = 2 and shot = 3): Multiple Choice Grade vs. Effective Model Size.
Bottom row (shot = 2 and shot = 3): - Brier Score vs. Effective Model Size.

task
- irony_identification
- logical_args
- sports_understanding
- swahili_english_proverbs

Part I:    Intuition on emergent abilities

Part II:    Empirically prove hypothesis InstructGPT/GPT-3

Part III:   Meta-analysis of emergent abilities

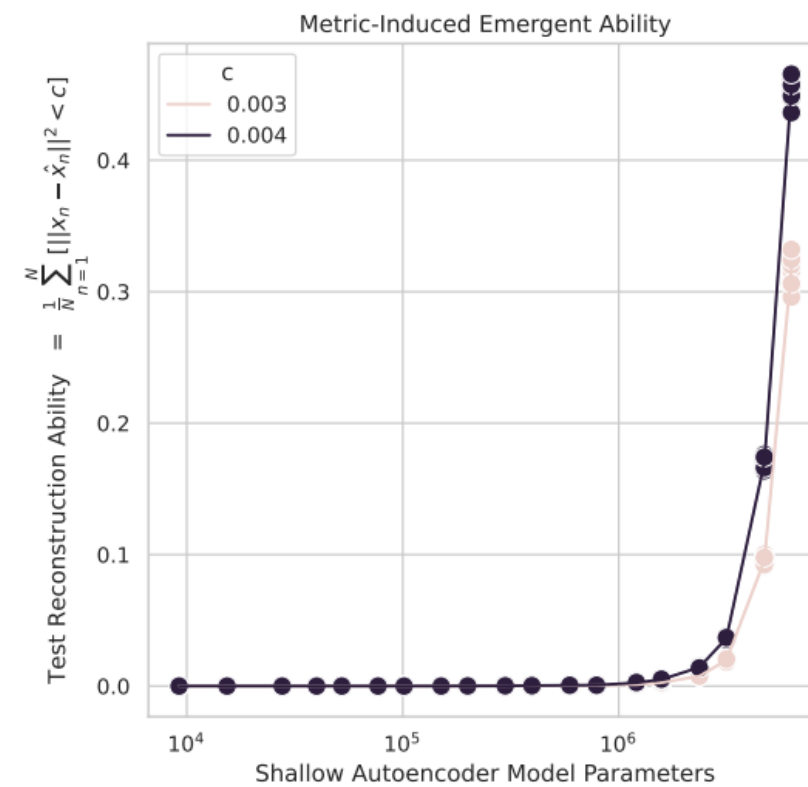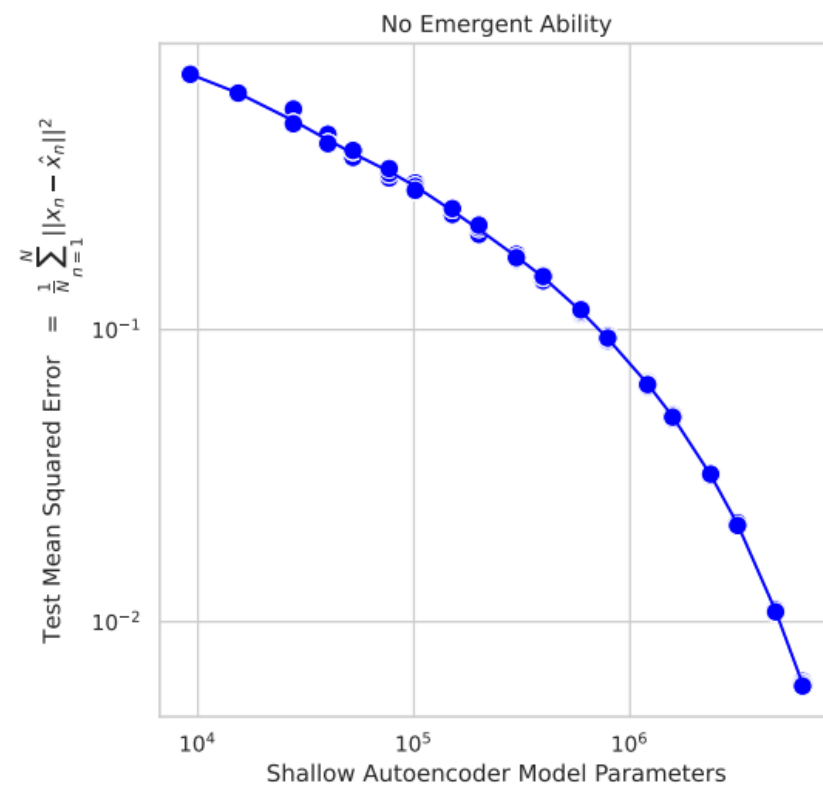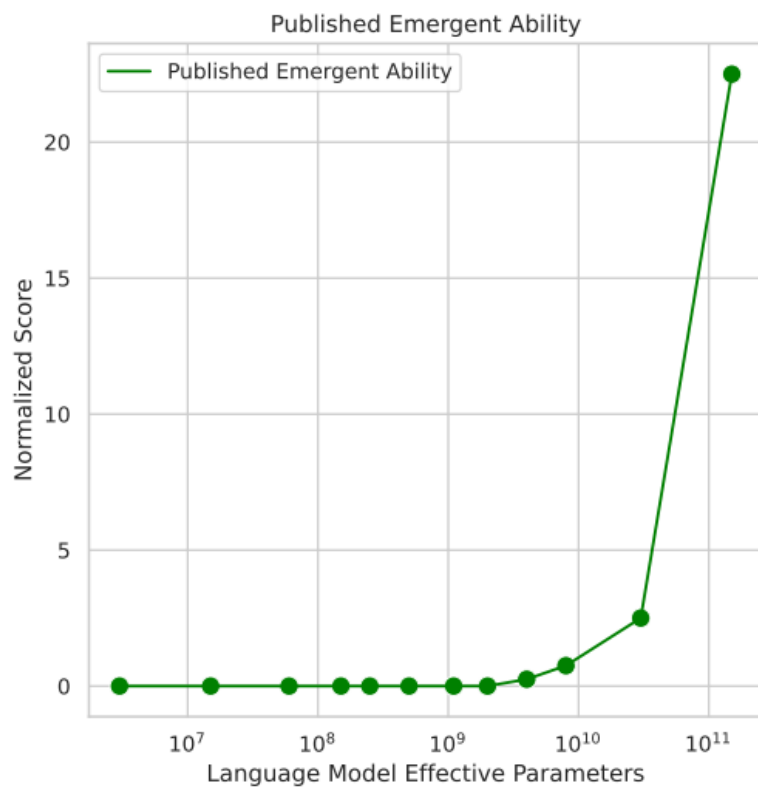Part IV:   Induce emergent abilities on vision models

Part I:     Intuition on emergent abilities

Part II:    Empirically prove hypothesis InstructGPT/GPT-3
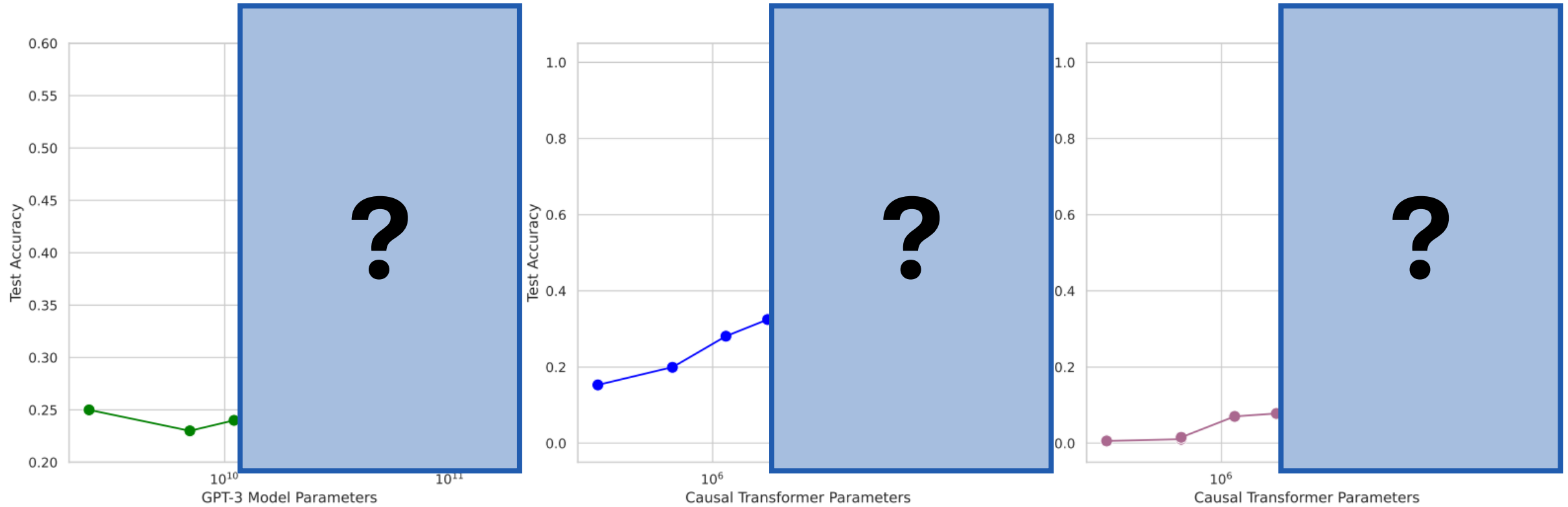
Part III:   Meta-analysis of emergent abilities

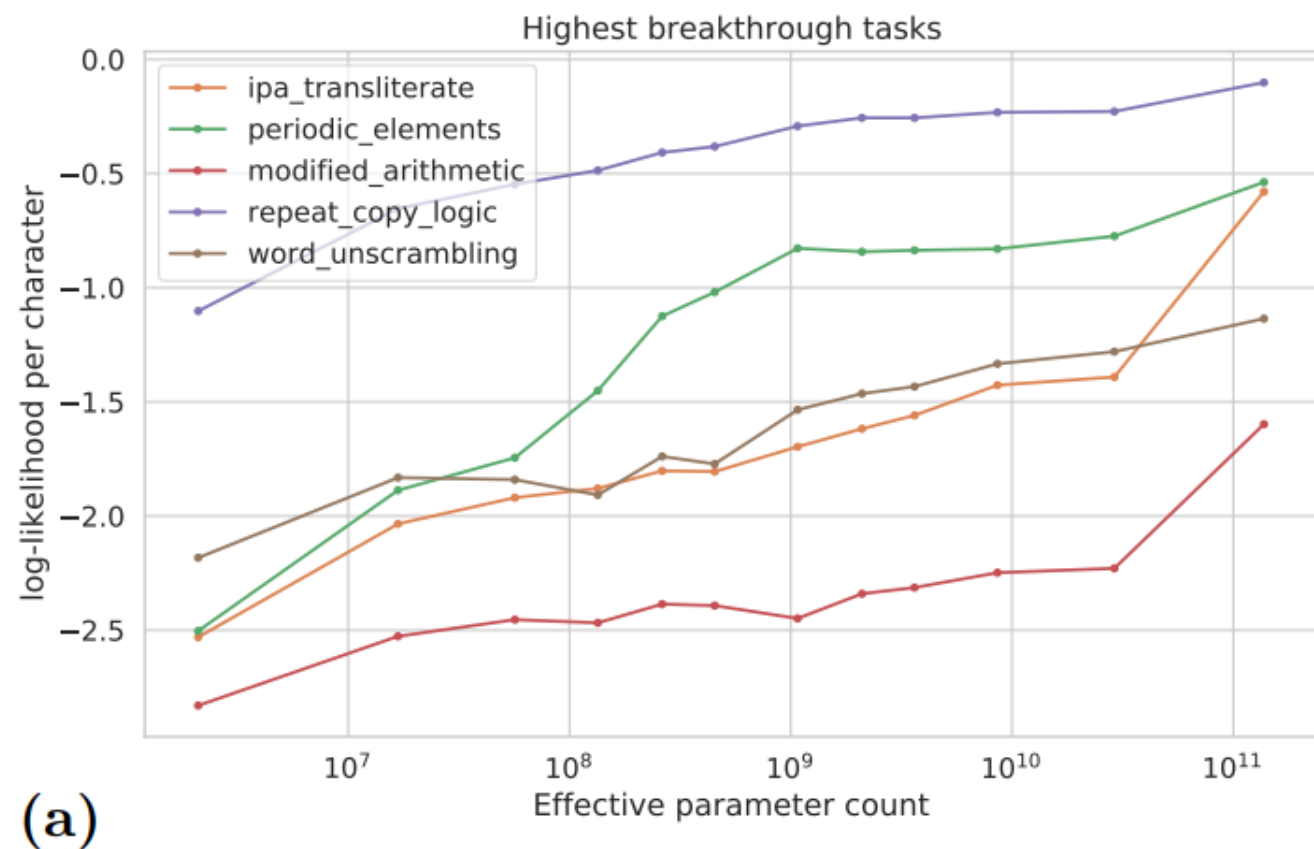Part IV:   Induce emergent abilities on vision models

# Shallow non-linear autoencoder for CIFAR100

# Transformer for classifying Omniglot characters

# BIG Benchmark

# GPT-4 Technical Report



**Capability prediction on 23 coding problems**

# Are Emergent Abilities a Mirage?

Emergent abilities only occur with certain metrics.

Those metrics are the ones that matter.

# Are Emergent Abilities a Mirage?



Plots have large jumps due
to log-scaled x-axis.

Linear scaling of x-axis
also shows jumps.

# Are Emergent Abilities a Mirage?



X-axis is not sampled densely enough.

The trend cannot be extrapolated.

# Discussion

**Strengths**:
- Multiple arguments to support their hypothesis.
- Clear explanation for unpredictable trends.

**Weaknesses**:
- Unpredictability of improvement.

# Are Emergent Abilities of Large Language Models a Mirage?

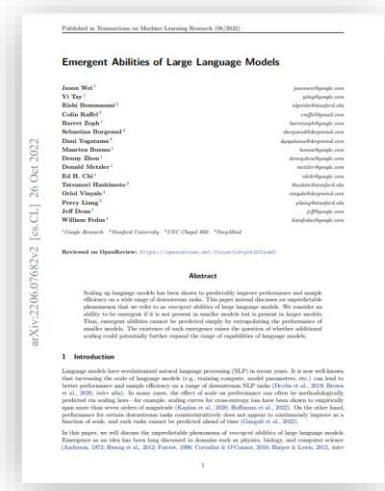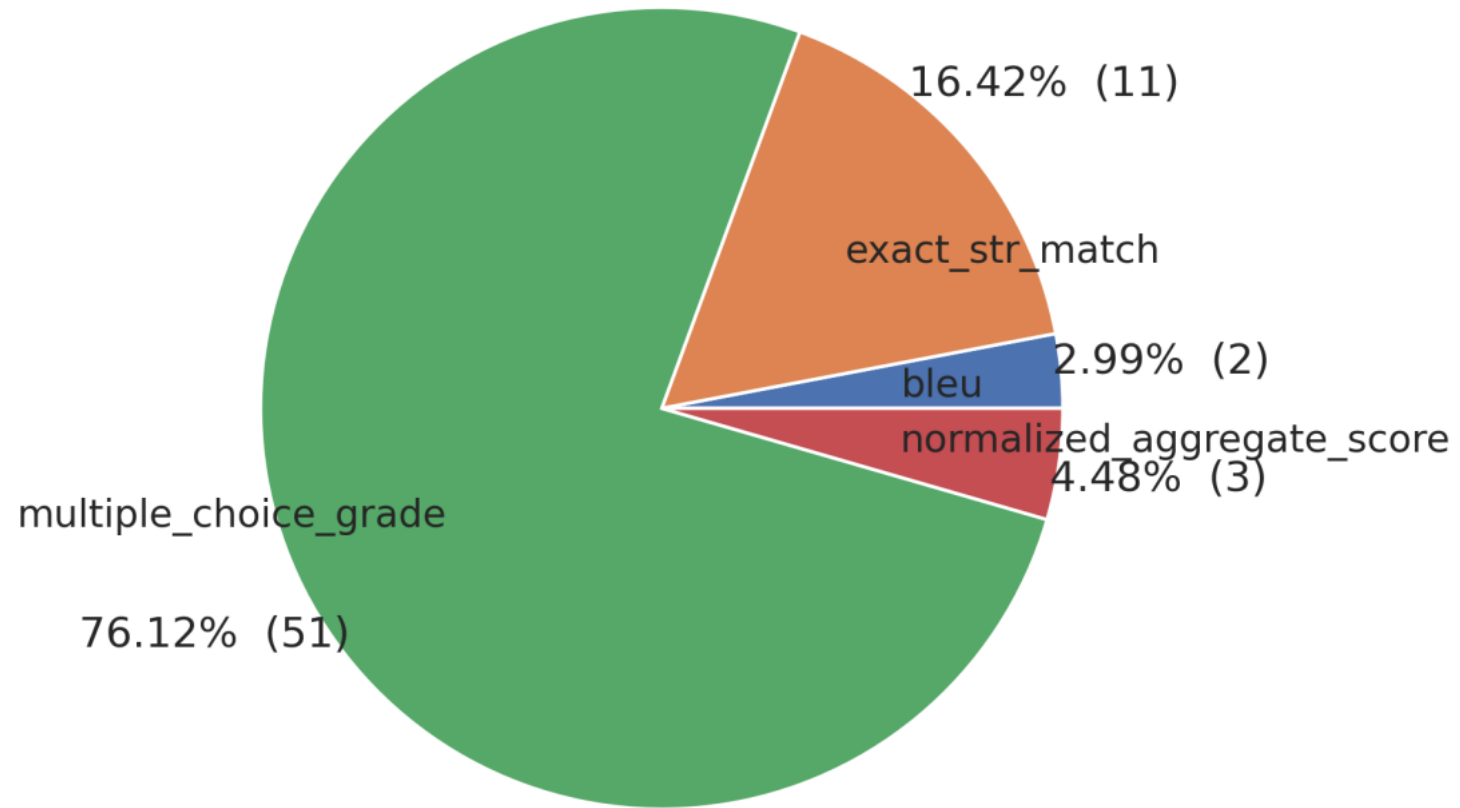Rylan Schaeffer, Brando Miranda, and Sanmi Koyejo

Computer Science, Stanford University

## Abstract

Recent work claims that large language models display *emergent abilities*, abilities not present in smaller-scale models that are present in larger-scale models. What makes emergent abilities intriguing is two-fold: their *sharpness*, transitioning seemingly instantaneously from not present to present, and their *unpredictability*, appearing at seemingly unforeseeable model scales. Here, we present an alternative explanation for emergent abilities: that for a particular task and model family, when analyzing fixed model outputs, emergent abilities appear due the researcher's choice of metric rather than due to fundamental changes in model behavior with scale. Specifically, nonlinear or discontinuous metrics produce apparent emergent abilities, whereas linear or continuous metrics produce smooth, continuous, predictable changes in model performance. We present our alternative explanation in a simple mathematical model, then test it in three complementary ways: we (1) make, test and confirm three predictions on the effect of metric choice using the InstructGPT/GPT-3 family on tasks with claimed emergent abilities, (2) make, test and confirm two predictions about metric choices in a meta-analysis of emergent abilities on BIG-Bench; and (3) show how to choose metrics to produce never-before-seen seemingly emergent abilities in multiple vision tasks across diverse deep networks. Via all three analyses, we provide evidence that alleged emergent abilities evaporate with different metrics or with better statistics, and may not be a fundamental property of scaling AI models.

# Are Emergent Abilities of LLMs a Mirage

By Rylan Schaeffer, Brando Miranda, and Sanmi Koyejo

Presented by Chao Chen (Michelle)
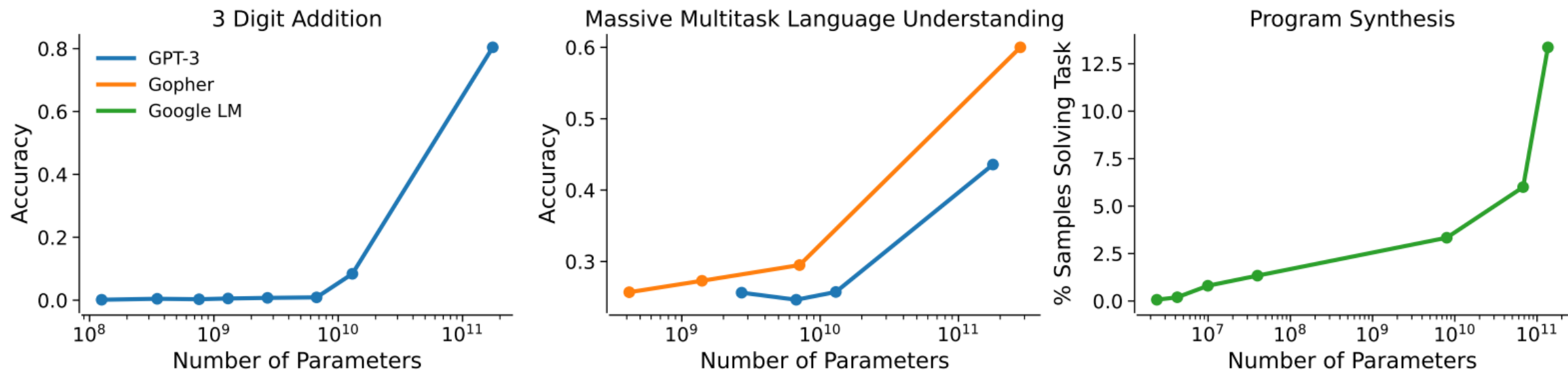
35

# Predictability and Surprise in LGM



**Fig. 2** Three examples of abrupt specific capability scaling described in Section 2.2, based on three different models: GPT-3 (blue), Gopher (orange), and a Google language model (green). **(Left)** 3-Digit addition with GPT-3 [11]. **(Middle)** Language understanding with GPT-3 and Gopher [62]. **(Right)** Program synthesis with Google language models [4].