# Building Language Models means making decisions



MMLU

- Dataset A
- Dataset B
- × Final checkpoint

Accuracy

Training Step

**Data decisions:** Train a **pair** models and compare the scores

MMLU

RC Accuracy

error 3.7%

*E.g. a 13B trained to 5T tokens*

Compute (FLOPs)

- ● Scaling Law Models
- × Predicted 13B Model
- ● Real 13B Model
- — Scaling Law Fit

**Scaling laws:** Train many **small** models and extrapolate the performance

DataDecide: How to Predict Best Pretraining Data with Small Experiments (ICML, 2025)
Establishing Task Scaling Laws via Compute-Efficient Model Ladders (COLM, 2025)

✦Ai2

# Downstream tasks are now core to building models …
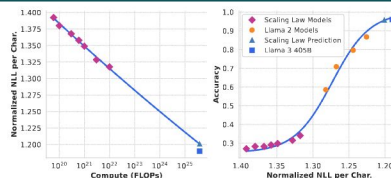
## … some tasks are useful!



**Figure 4  Scaling law forecast for ARC Challenge.** *Left:* Normalized negative log-likelihood of the correct answer on the ARC Challenge benchmark as a function of pre-training FLOPs. *Right:* ARC Challenge benchmark accuracy as a function of the normalized negative log-likelihood of the correct answer. This analysis enables us to predict model performance on the ARC Challenge benchmark before pre-training commences. See text for details.
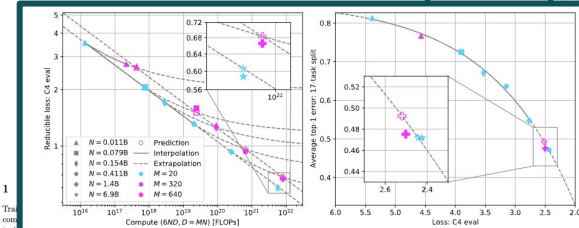
## … some tasks aren't predictable



Figure 1: **Reliable scaling with over-training and on downstream error prediction.** *(left)* We fit a scaling law for model validation loss, parameterized by (i) a token multiplier $M = N/D$, which is the ratio of training tokens $D$ to parameters $N$ and (ii) the compute $C$ in FLOPs used to train a model, approximated by $C = 6ND$. Larger values of $M$ specify more over-training. We are able to extrapolate, in both $N$ and $M$, the validation performance of models requiring more than $300\times$ the training compute used to construct the scaling law. *(right)* We also fit a scaling law to predict average downstream top-1 error as a function of validation loss. We find that fitting scaling laws for downstream error benefits from using more expensive models when compared to fitting for loss prediction. We predict the average error over 17 downstream tasks for models trained with over $20\times$ the compute. For this figure, we train all models on RedPajama [112].

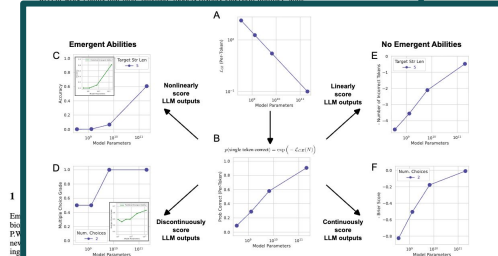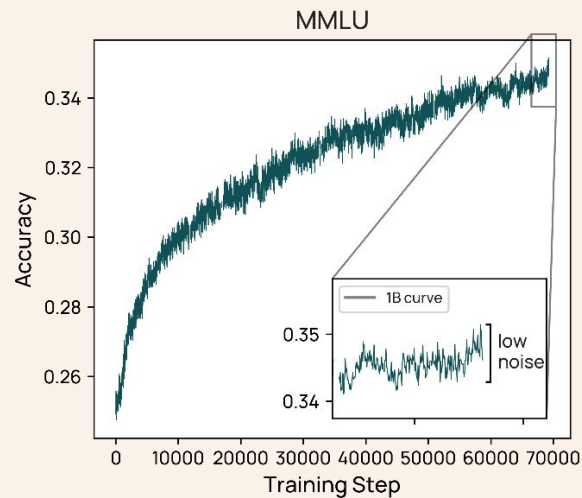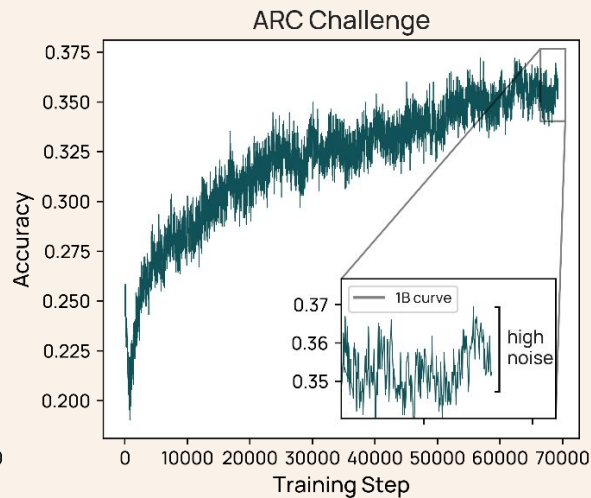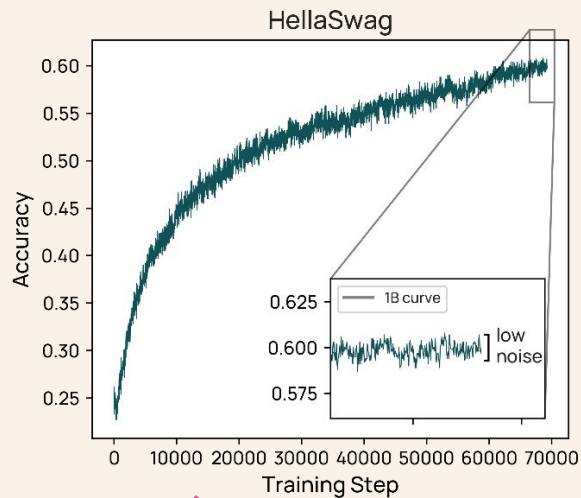## … some metrics hide real capability



Figure 2: **Emergent abilities of large language models are created by the researcher's chosen metrics, not unpredictable changes in model behavior with scale.** (A) Suppose the per-token cross-entropy loss decreases monotonically with model scale, e.g., $\mathcal{L}_{CE}$ scales as a power law. (B) The per-token probability of selecting the correct token asymptotes towards 1. (C) If the researcher scores models' outputs using a nonlinear metric such as Accuracy (which requires a sequence of tokens to *all* be correct), the metric choice nonlinearly scales performance, causing performance to change sharply and unpredictably in a manner that qualitatively matches published emergent abilities (inset). (D) If the researcher instead scores models' outputs using a discontinuous metric such as Multiple Choice Grade (akin to a step function), the metric choice discontinuously scales performance, again causing performance to change sharply and unpredictably. (E) Changing from a nonlinear metric to a linear metric such as Token Edit Distance, scaling shows smooth, continuous and predictable improvements, ablating the emergent ability. (F) Changing from a discontinuous metric to a continuous metric such as Brier Score again reveals smooth, continuous and predictable improvements in task performance. Consequently, emergent abilities are created by the researcher's choice of metrics, not fundamental changes in model family behavior on specific tasks with scale.

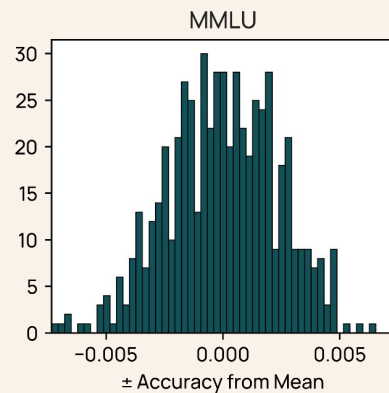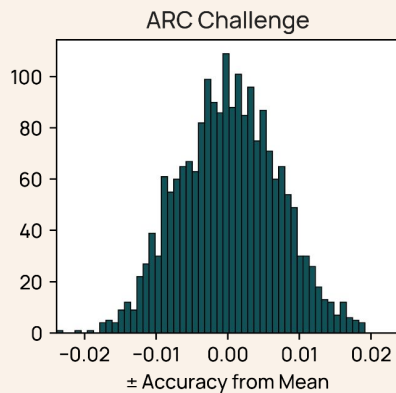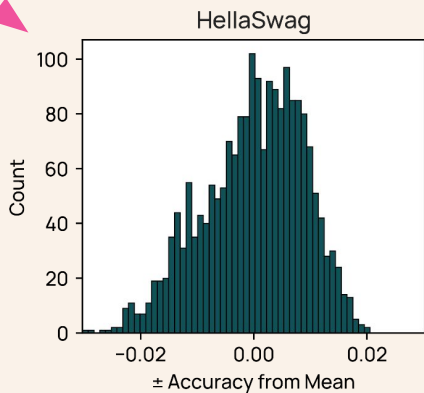The Llama 3 Herd of Models (preprint, 2024)
Language models scale reliably with over-training and on downstream tasks (ICLR, 2025)
Are Emergent Abilities of Large Language Models a Mirage? (NeurIPS, 2023)
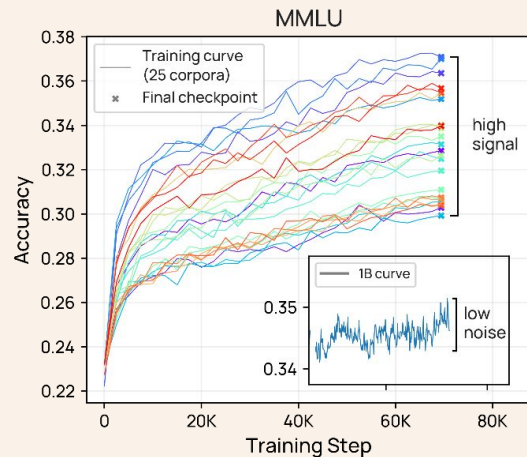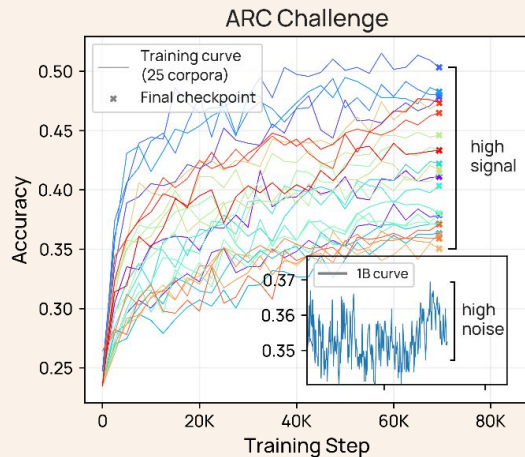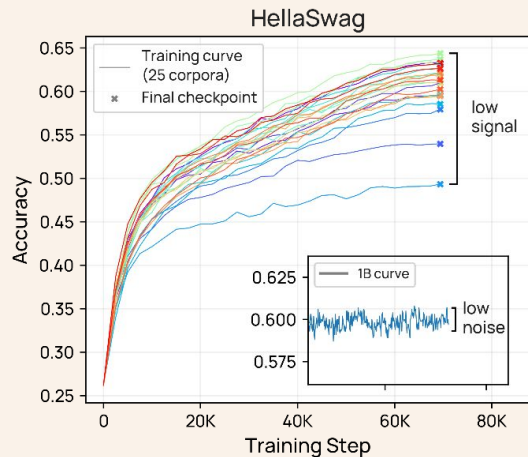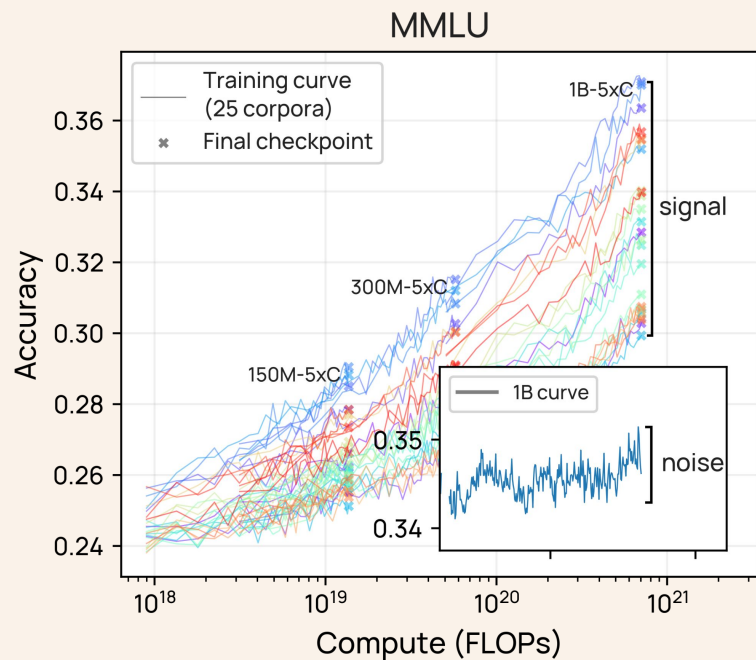
❖Ai2

# Why do so many predictions fail - but some don't?

HellaSwag — ARC Challenge — MMLU

Final 20% of checkpoints

❖Ai2

# ... but inter-checkpoint variance is not the whole story! We need to measure both signal and noise



Ai2

**signal:**

$$\text{Rel. Dispersion}(M) = \max_{j,k} |m_j - m_k|/\bar{m}$$

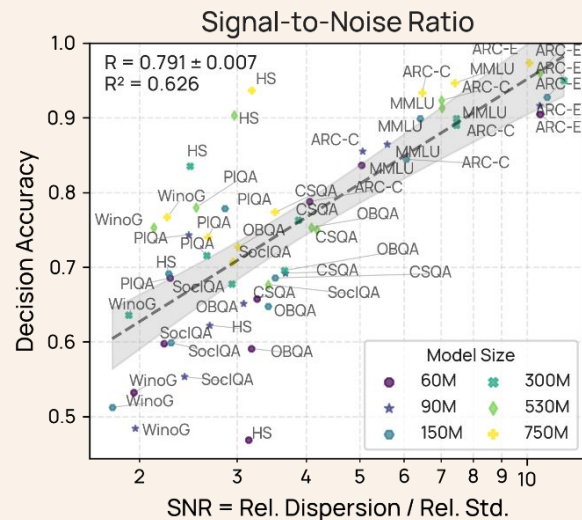**noise:**

$$\text{Rel. Std.}(m) = \sqrt{\frac{1}{n-1} \sum_{i=1}^{n} (m_i - \bar{m})^2}/\bar{m}$$

$$\text{Signal-to-Noise Ratio} =$$

$$\frac{\text{Rel. Dispersion(final train checkpoint)}}{\text{Rel. Std.(final } n \text{ train checkpoints)}}$$

# Only signal or noise alone do not explain rank agreement from small to large scale... ... but the signal-to-noise ratio does!

# Predicting task performance using
# scaling laws is sensitive to noise!

# Predicting task performance using scaling laws is sensitive to noise!

Scatter plot with x-axis "Noise (final 30 train checkpoints)" and y-axis "Scaling Law Prediction Error". Annotations on plot:
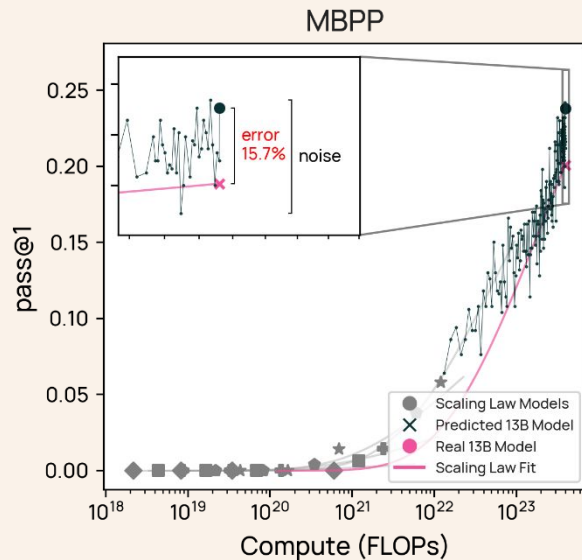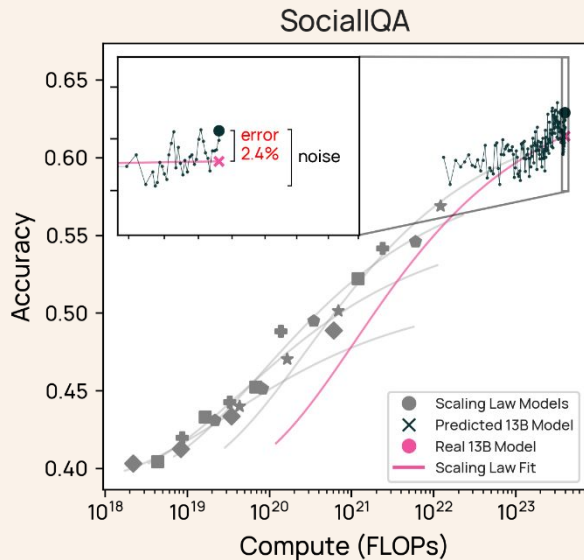
R = 0.653 ± 0.068
R² = 0.426

Labeled points: HellaSwag, Jeopardy, All Tasks, SocialIQA, Knowledge Tasks, MMLU, TriviaQA, MBPP+, Code Tasks, Math Tasks, MedMCQA, Minerva MATH.

Model Size: 13B

# We can use the signal-to-noise ratio to improve our benchmarks

➤ **intervention 1: filter subtasks with high SNR**

➤ intervention 2: smooth intermediate checkpoints

➤ intervention 3: select metrics with high SNR

### MMLU Signal-to-noise Ratio



Highest signal-to-noise ratio is top 16 MMLU subtasks

Included MMLU Subtask

— Subtasks sorted by SNR     — Subtasks sorted randomly

❖ Ai2

# We can use the signal-to-noise ratio to improve our benchmarks

▷ intervention 1: filter subtasks with high SNR

► **intervention 2: smooth intermediate checkpoints**

▷ intervention 3: select metrics with high SNR



❖ Ai2

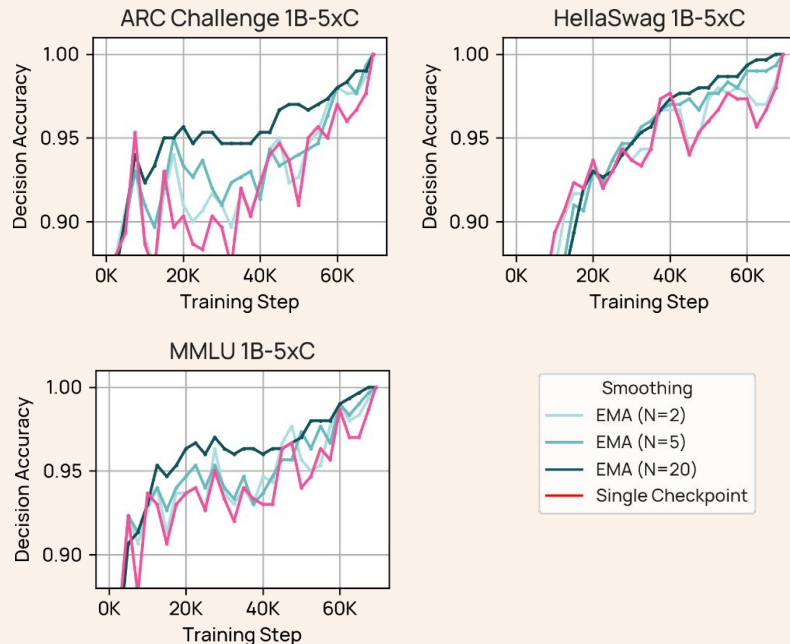# We can use the signal-to-noise ratio to improve our benchmarks

➤ intervention 1: filter subtasks with high SNR

➤ intervention 2: smooth intermediate checkpoints

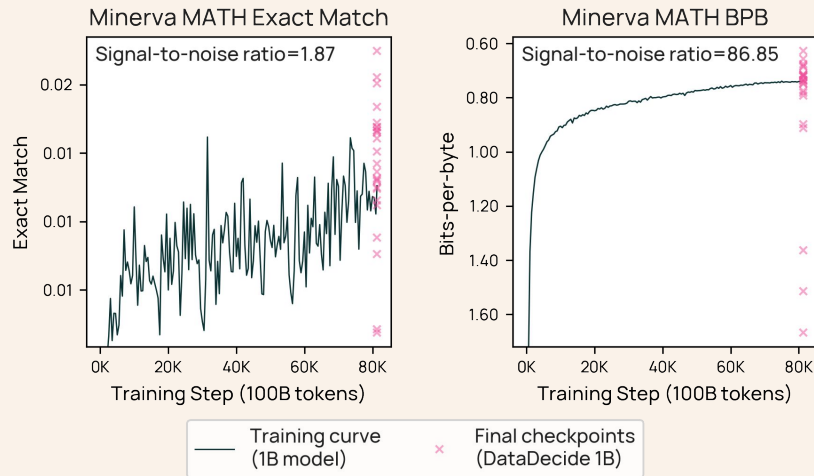➤ **intervention 3: select metrics with high SNR**



Minerva MATH Exact Match
Signal-to-noise ratio=1.87

Minerva MATH BPB
Signal-to-noise ratio=86.85

Training curve (1B model) — Final checkpoints (DataDecide 1B)

✦ Ai2

# Thank you!

Learn more at our poster:
### Wednesday, Dec 3 at 4:30 in Hall C,D,E.

**Contact:** davidh@allenai.org