# What's so interesting about models with recurrent depth?

**Jonas Geiping**

# Scaling up Test-Time Compute with Latent Reasoning:
# A Recurrent Depth Approach

Jonas Geiping [1]  Sean McLeish [2]  Neel Jain [2]  John Kirchenbauer [2]  Siddharth Singh [2]  Brian R. Bartoldson [3]

Bhavya Kailkhura [3]  Abhinav Bhatele [2]  Tom Goldstein [2]

Paper:                           https://www.arxiv.org/abs/2502.05171

Code:                           github.com/seal-rg/recurrentpretraining

Model:                          huggingface.co/tomg-group-umd/huginn-0125

Data (still uploading):   huggingface.co/datasets/tomg-group-umd/huginn-dataset

# Key Concepts

# What is Test-Time Compute?

Improving the performance of a machine learning model by using more computation at test-time/ for inference.

# What is "Reasoning"

Learning to pattern match abstract concepts such as meta-strategies, logic, abstraction and plans, instead of just memorizing facts.

# What is *Verbalized* Reasoning

Strategized optimal cargo placement across aircraft zones.                    ⏱ 2m, 9s

To solve this problem, I need to find a valid arrangement of the five cargo items (A, B, C, D, E) across the three zones (Front, Middle, Rear) that satisfies all of the given constraints.

Let me first summarize the constraints and items:

**Aircraft Zones:**

- Front: max 2000kg, 8m³
- Middle: max 3000kg, 12m³
- Rear: max 2500kg, 10m³

**Cargo Items:**

- Item A: 1200kg, 5m³, hazard class H1, priority 2, CG impact +5 per 100kg
- Item B: 800kg, 4m³, hazard class H2, priority 1, CG impact -3 per 100kg
- Item C: 1500kg, 6m³, hazard class H1, priority 3, CG impact +2 per 100kg

# *Verbalized* Reasoning

Strategized optimal cargo placement across aircraft zones.

⏱ 2m, 9s

To solve this problem, I need to find a valid arrangement of the five cargo items (A, E
across the three zones (Front, Middle, Rear) that satisfies all of the given constraints.

Let me first summarize the constraints and items:

**Aircraft Zones:**

- Front: max 2000kg, 8m³
- Middle: max 3000kg, 12m³
- Rear: max 2500kg, 10m³

**Cargo Items:**

- Item A: 1200kg, 5m³, hazard class H1, priority 2, CG impact +5 per 100kg
- Item B: 800kg, 4m³, hazard class H2, priority 1, CG impact -3 per 100kg
- Item C: 1500kg, 6m³, hazard class H1, priority 3, CG impact +2 per 100kg

This arrangement doesn't work. Let's try another.

Attempt 2:

- B and E in the Rear (priority 1)
- A and D in the Front (ensuring H1 and H3 aren't adjacent)
- C in the Middle

Let's check:

**Front:** Items A + D (1200kg + 1800kg = 3000kg, 5m³ + 7m³ = 12m³) - exceeds weight limit (3000kg > 2000kg)

This doesn't work either. Let's try another arrangement.

Attempt 3:

- B and E in the Rear (priority 1)
- A in the Front

**4** ___

# *Verbalized* Reasoning

Strategized optimal cargo placement across aircraft zones.　🕐 2m, 9s

To solve this problem, I need to find a valid arrangement of the five cargo items (A, B,
across the three zones (Front, Middle, Rear) that satisfies all of the given constraints.

Let me first summarize the constraints and items:

**Aircraft Zones:**

- Front: max 2000kg, 8m³
- Middle: max 3000kg, 12m³
- Rear: max 2500kg, 10

**Cargo Items:**

- Item A: 1200kg, 5m³,
- Item B: 800kg, 4m³, h
- Item C: 1500kg, 6m³, l

Strategized optimal cargo placement across aircraft zones.　🕐 2m, 9s

This arrangement doesn't work. Let's try another.

Attempt 2:

- B and E in the Rear (priority 1)
- A and D in the Front (ensuring H1 and H3 aren't adjacent)

Strategized optimal cargo placement across aircraft zones.　🕐 2m, 9s

This gives us these potential arrangements when B + A are in the Rear:　　　　7m³ = 12m³) - exceeds weight limit (3000kg >

1. Front: E, Middle: C+D (C is H1, D is H3 - invalid)
2. Front: C, Middle: E+D (C is H1, D is H3 - invalid)　　　　　　　　　　　nt.
3. Front: D, Middle: E+C (A is H1 in Rear, D is H3 in Front, C is H1 in Middle - both A and C are not adjacent to D, so this is valid!)
4. Front: E+C, Middle: D (C is H1, D is H3 - invalid)

So, we have a valid arrangement!

- Front: D (1800kg, 7m³, H3)
- Middle: E (900kg, 3m³, H2, P1) + C (1500kg, 6m³, H1) = 2400kg, 9m³
- Rear: B (800kg, 4m³, H2, P1) + A (1200kg, 5m³, H1) = 2000kg, 9m³

This satisfies:

- Weight and volume constraints for all zones

4 —

# *Verbalized* Reasoning

To solve this problem, I need to find a valid arrangement of the five cargo items (A, B...
across the three zones (Front, Middle, Rear) that satisfies all of the given constraints.

Let me first summarize the constraints and items:

**Aircraft Zones:**

- Front: max 2000kg, 8m³
- Middle: max 3000kg, 12m³
- Rear: max 2500kg, 10...

**Cargo Items:**

- Item A: 1200kg, 5m³, ...
- Item B: 800kg, 4m³, h...
- Item C: 1500kg, 6m³, ...

Strategized optimal cargo placement across aircraft zones.

This gives us these potential arrangements when B + A are i...

1. Front: E, Middle: C+D (C is H1, D is H3 - invalid)
2. Front: C, Middle: E+D (C is H1, D is H3 - invalid)
3. Front: D, Middle: E+C (A is H1 in Rear, D is H3 in Front, adjacent to D, so this is valid!)
4. Front: E+C, Middle: D (C is H1, D is H3 - invalid)

So, we have a valid arrangement!

- Front: D (1800kg, 7m³, H3)
- Middle: E (900kg, 3m³, H2, P1) + C (1500kg, 6m³, H1) = 2...
- Rear: B (800kg, 4m³, H2, P1) + A (1200kg, 5m³, H1) = 200...

This satisfies:

- Weight and volume constraints for all zones

second (from the Middle). So in terms of priority item accessibility, all three arrangements are equally optimal.

Since all three arrangements satisfy all constraints and have the same CG impact, and are equally optimal in terms of priority item accessibility, any of them could be considered the "optimal" arrangement. For the sake of providing a specific answer, I'll go with Arrangement 3:

- Front: D (1800kg, 7m³, H3)
- Middle: B (800kg, 4m³, H2, P1) + A (1200kg, 5m³, H1)
- Rear: E (900kg, 3m³, H2, P1) + C (1500kg, 6m³, H1)

So the optimal arrangement of the five cargo items across the three zones is:

- Front Zone: Item D
- Middle Zone: Items B and A
- Rear Zone: Items E and C

To solve the Automated Aircraft Loading Problem, I need to find an arrangement that satisfies all constraints while optimizing for priority access.

Here's the optimal solution:

**Front Zone:**

- Item D (1800kg, 7m³, hazard class H3)

**Middle Zone:**

- Item B (800kg, 4m³, hazard class H2, priority 1)
- Item A (1200kg, 5m³, hazard class H1)

**Rear Zone:**

- Item E (900kg, 3m³, hazard class H2, priority 1)
- Item C (1500kg, 6m³, hazard class H1)

This solution satisfies all constraints:

# *Continuous* Reasoning?

Moving the reasoning chain into the model's representation space.

# What is recurrent depth?

A model that recurs "in depth" re-uses its layers or other subcomponents to create a deeper or shallower model.

# What is recurrent depth?

A model that recurs "in depth" re-uses its layers or other subcomponents to create a deeper or shallower model.

Universal Transformers

# What is recurrent depth?

A model that recurs "in depth" re-uses its layers or other subcomponents to create a deeper or shallower model.

Universal Transformers

Looped Transformers

# What is recurrent depth?

A model that recurs "in depth" re-uses its layers or other subcomponents to create a deeper or shallower model.

Universal Transformers

Representation Recycling

Looped Transformers

# What is recurrent depth?

A model that recurs "in depth" re-uses its layers or other subcomponents to create a deeper or shallower model.

Universal Transformers

RNNs with Adaptive Computation Time

Representation Recycling

Looped Transformers

# What is recurrent depth?

A model that recurs "in depth" re-uses its layers or other subcomponents to create a deeper or shallower model.

Universal Transformers

RNNs with Adaptive Computation Time

Equilibrium Models

Representation Recycling

Looped Transformers

# What is recurrent depth?

A model that recurs "in depth" re-uses its layers or other subcomponents to create a deeper or shallower model.

Universal Transformers

RNNs with Adaptive Computation Time

Equilibrium Models

Representation Recycling

Diffusion Models

Looped Transformers

# What is recurrent depth?

A model that recurs "in depth" re-uses its layers or other subcomponents to create a deeper or shallower model.

Universal Transformers

RNNs with Adaptive Computation Time

Equilibrium Models

Representation Recycling

Diffusion Models

Looped Transformers

Iterative Refinement

# What is recurrent depth?

A model that recurs "in depth" re-uses its layers or other subcomponents to create a deeper or shallower model.

Universal Transformers          RNNs with Adaptive Computation Time

Equilibrium Models                              Representation Recycling

Diffusion Models                Looped Transformers
Iterative Refinement                      Implicit Neural nets

# What is recurrent depth?

A model that recurs "in depth" re-uses its layers or other subcomponents to create a deeper or shallower model.

Universal Transformers          RNNs with Adaptive Computation Time

**Hopfield-Amari Models?**

Equilibrium Models                                                          Representation Recycling

Diffusion Models                        Looped Transformers
                                                                    Implicit Neural nets
Iterative Refinement

# Why use recurrent-depth as a framework for test-time compute?

- No supervision on intermediate steps, so no CoT data needed with the right training objective
- No increased context length, linear complexity increase
- Recurrent-depth models have less parameters
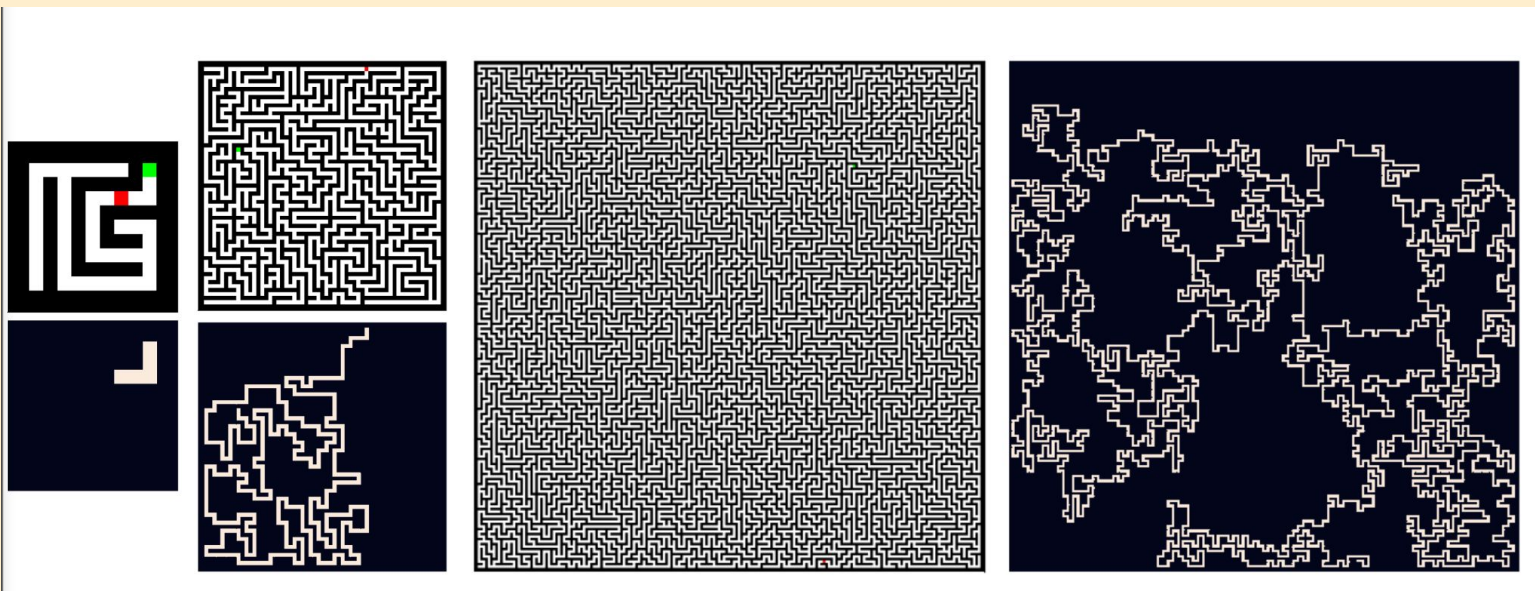
- Recurrent-depth models are *compute-heavy*

# A good prior for "Reasoning"?

- Easy to learn iterative *algorithms*
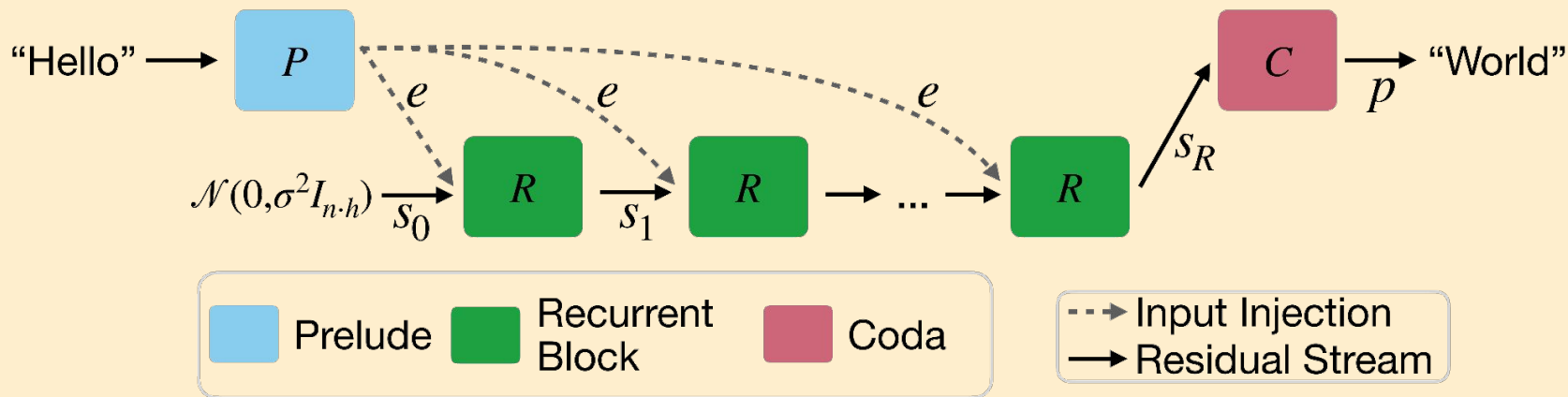- Harder for the model to memorize

# A good prior for "Reasoning"?

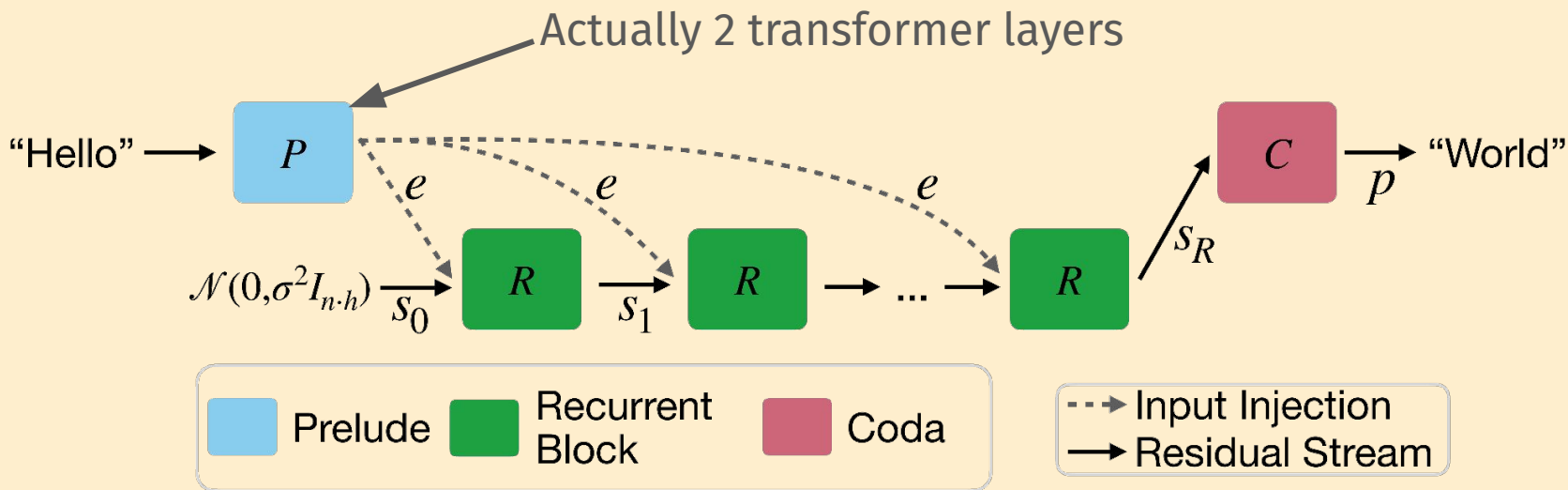- Easy to learn iterative *algorithms*
- Harder for the model to memorize

# A good prior for "Reasoning"?

- Easy to learn iterative *algorithms*
- Harder for the model to memorize

Bansal*, Schwarzschild*, Borgnia, Emam, Huang, Goldblum, Goldstein
*End-to-end Algorithm Synthesis with Recurrent Networks: Logical Extrapolation Without Overthinking*

# A good prior for "Reasoning"?

- Easy to learn iterative *algorithms*
- Harder for the model to memorize

9

# A scalable recurrent (depth) architecture

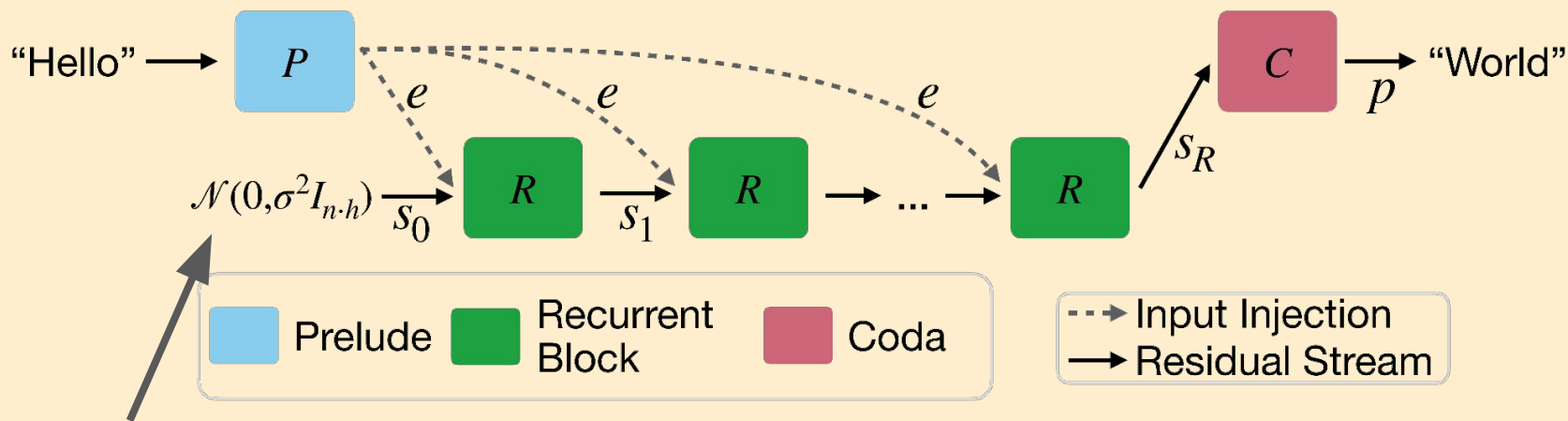# Why did we call this *latent* recurrent depth?

# Why did we call this *latent* recurrent depth?



Actually 2 transformer layers

"Hello" → $P$

$\mathcal{N}(0, \sigma^2 I_{n \cdot h})$ →$s_0$→ $R$ →$s_1$→ $R$ → ... → $R$ →$s_R$→ $C$ →$p$→ "World"

$e$   $e$   $e$

Prelude   Recurrent Block   Coda

- - → Input Injection
→ Residual Stream

# Why did we call this *latent* recurrent depth?



The "input injection" is crucial for stability

# Why did we call this *latent* recurrent depth?



"Hello" → $P$

$e$   $e$   $e$

$\mathcal{N}(0, \sigma^2 I_{n \cdot h})$   $s_0$ → $R$ → $s_1$ → $R$ → ... → $R$   $s_R$ → $C$   $p$ → "World"

Prelude   Recurrent Block   Coda

- - → Input Injection
⟶ Residual Stream

Diffusion Model connection,
appears based on findings for
path independence

# What does that mean?

# Training Objective



$$\mathcal{L}(\theta) = \mathbb{E}_{\mathbf{x} \in X} \mathbb{E}_{r \sim \Lambda} L\left(m_\theta(\mathbf{x}, r), \mathbf{x}'\right)$$

For every training sequence
- Sample a random* number of steps r to recur
- Compute r forward steps
- Compute loss based on the last k steps.

# Training Objective Alternatives?

- **Universal Transformers:** Train with halting module
- **Equilibrium Models:** Iterate to convergence, differentiate fixed point based on IFT
- **Weight-shared models:** Just fix number of steps, train as normal
- **Diffusion Model:** Train to denoise target hidden states

# Actually Training a Model at Scale

# How do we show that this *actually* scales?

- **100m parameter, 10B tokens prototypes all work great...**

# How do we show that this *actually* scales?

- 100m parameter, 10B tokens prototypes all work great...

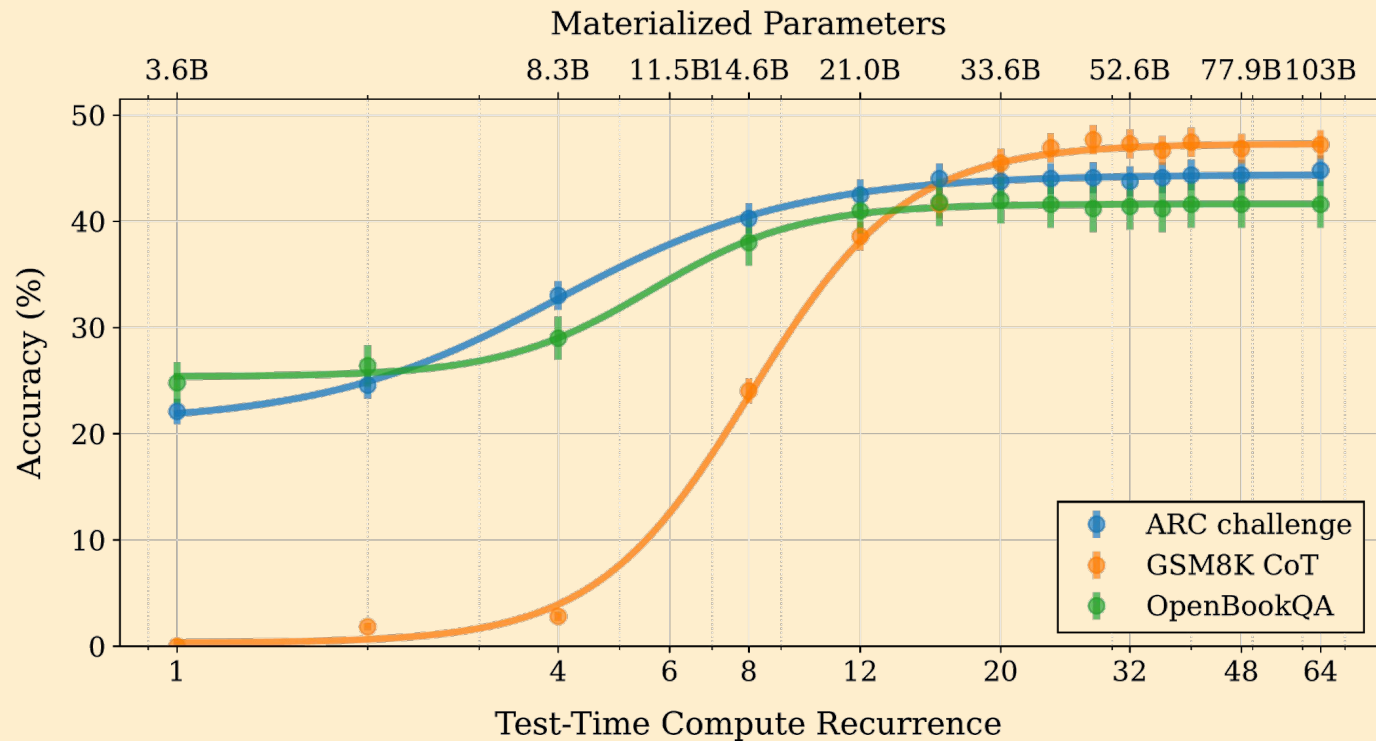To show that this worked we define a more convincing

 target,  *Huginn-3.5B*

# How do we show that this *actually* scales?

- 100m parameter, 10B tokens prototypes all work great...

To show that this worked we define a more convincing

target,  *Huginn-3.5B*

- 2 + 4 + 2 layers, tied embeds, 3.5B parameters
- We target 1T tokens of a wide pretraining mix
- Will it actually train to be a (semi)-competitive language and reasoning model?
- Who will give us compute ...

# Results

Scaling up Test-Time Compute with Recurrent Depth

# Standard benchmarks

| Model | Param | Tokens | ARC-E | ARC-C | HellaSwag | MMLU | OBQA | PiQA | SciQ | WinoGrande |
|---|---|---|---|---|---|---|---|---|---|---|
| random | | | 25.0 | 25.0 | 25.0 | 25.0 | 25.0 | 50.0 | 25.0 | 50.0 |
| Amber | 7B | 1.2T | 65.70 | 37.20 | 72.54 | 26.77 | 41.00 | 78.73 | 88.50 | 63.22 |
| Pythia-2.8b | 2.8B | 0.3T | 58.00 | 32.51 | 59.17 | 25.05 | 35.40 | 73.29 | 83.60 | 57.85 |
| Pythia-6.9b | 6.9B | 0.3T | 60.48 | 34.64 | 63.32 | 25.74 | 37.20 | 75.79 | 82.90 | 61.40 |
| Pythia-12b | 12B | 0.3T | 63.22 | 34.64 | 66.72 | 24.01 | 35.40 | 75.84 | 84.40 | 63.06 |
| OLMo-1B | 1B | 3T | 57.28 | 30.72 | 63.00 | 24.33 | 36.40 | 75.24 | 78.70 | 59.19 |
| OLMo-7B | 7B | 2.5T | 68.81 | 40.27 | 75.52 | 28.39 | 42.20 | 80.03 | 88.50 | 67.09 |
| OLMo-7B-0424 | 7B | 2.05T | 75.13 | 45.05 | 77.24 | 47.46 | 41.60 | 80.09 | 96.00 | 68.19 |
| OLMo-7B-0724 | 7B | 2.75T | 74.28 | 43.43 | 77.76 | 50.18 | 41.60 | 80.69 | 95.70 | 67.17 |
| OLMo-2-1124 | 7B | 4T | 82.79 | 57.42 | 80.50 | 60.56 | 46.20 | 81.18 | 96.40 | 74.74 |
| Ours, ($r = 4$) | 3.5B | 0.8T | 49.07 | 27.99 | 43.46 | 23.39 | 28.20 | 64.96 | 80.00 | 55.24 |
| Ours, ($r = 8$) | 3.5B | 0.8T | 65.11 | 35.15 | 58.54 | 25.29 | 35.40 | 73.45 | 92.10 | 55.64 |
| Ours, ($r = 16$) | 3.5B | 0.8T | 69.49 | 37.71 | 64.67 | 31.25 | 37.60 | 75.79 | 93.90 | 57.77 |
| Ours, ($r = 32$) | 3.5B | 0.8T | 69.91 | 38.23 | 65.21 | 31.38 | 38.80 | 76.22 | 93.50 | 59.43 |

# Reasoning (grade-school math)

| Model | GSM8K | GSM8k CoT | Minerva MATH | MathQA |
|---|---|---|---|---|
| Random | 0.00 | 0.00 | 0.00 | 20.00 |
| Amber | 3.94/4.32 | 3.34/5.16 | 1.94 | 25.26 |
| Pythia-2.8b | 1.59/2.12 | 1.90/2.81 | 1.96 | 24.52 |
| Pythia-6.9b | 2.05/2.43 | 2.81/2.88 | 1.38 | 25.96 |
| Pythia-12b | 3.49/4.62 | 3.34/4.62 | 2.56 | 25.80 |
| OLMo-1B | 1.82/2.27 | 1.59/2.58 | 1.60 | 23.38 |
| OLMo-7B | 4.02/4.09 | 6.07/7.28 | 2.12 | 25.26 |
| OLMo-7B-0424 | 27.07/27.29 | 26.23/26.23 | 5.56 | 28.48 |
| OLMo-7B-0724 | 28.66/28.73 | 28.89/28.89 | 5.62 | 27.84 |
| OLMo-2-1124-7B | 66.72/66.79 | 61.94/66.19 | 19.08 | 37.59 |
| Our w/o sys. prompt ($r = 32$) | 28.05/28.20 | 32.60/34.57 | 12.58 | 26.60 |
| Our w/ sys. prompt ($r = 32$) | 24.87/38.13 | 34.80/42.08 | 11.24 | 27.97 |

# Reasoning (grade-school math)

| Model | Tokens | ARC-E | ARC-C | HellaSwag | MMLU | OBQA | PiQA | SciQ | WinoGrande | GSM8K CoT |
|---|---|---|---|---|---|---|---|---|---|---|
| Fixed-Depth Baseline | 0.18T | 46.42 | 26.96 | 37.34 | 24.16 | 29.60 | 64.47 | 73.20 | 51.78 | 1.82/2.20 |
| Ours, early ckpt, $(r = 32)$ | 0.18T | 53.62 | 29.18 | 48.80 | 25.59 | 31.40 | 68.88 | 80.60 | 52.88 | 9.02/10.24 |
| Ours, early ckpt, $(r = 1)$ | 0.18T | 34.01 | 23.72 | 29.19 | 23.47 | 25.60 | 53.26 | 54.10 | 53.75 | 0.00/0.15 |
| Ours, $(r = 32)$ | 0.8T | 69.91 | 38.23 | 65.21 | 31.38 | 38.80 | 76.22 | 93.50 | 59.43 | 34.80/42.08 |
| Ours, $(r = 1)$ | 0.8T | 34.89 | 24.06 | 29.34 | 23.60 | 26.80 | 55.33 | 47.10 | 49.41 | 0.00/0.00 |

# Scaling Test-Time Compute    vs    Scaling Pretraining

# What is the model doing?

Latent State Convergence $||x - x^*||$

# Token Trajectories

# Emergent Terminal Behaviors



Token: " deeper"

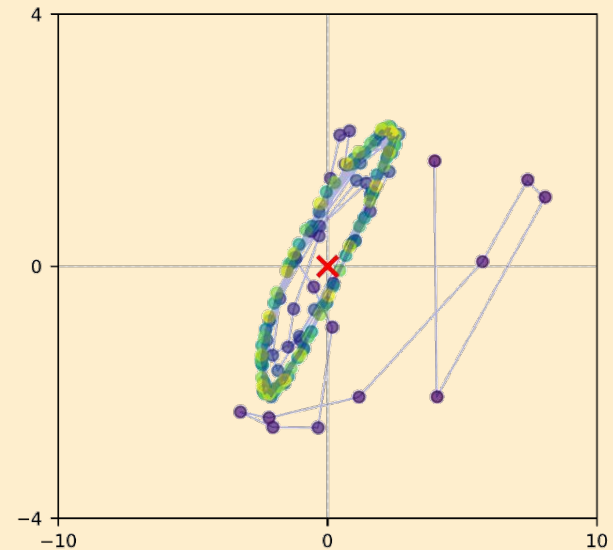# Emergent Terminal Behaviors



Token: " 3"
PC1-PC2   PC3-PC4   PC5-PC6
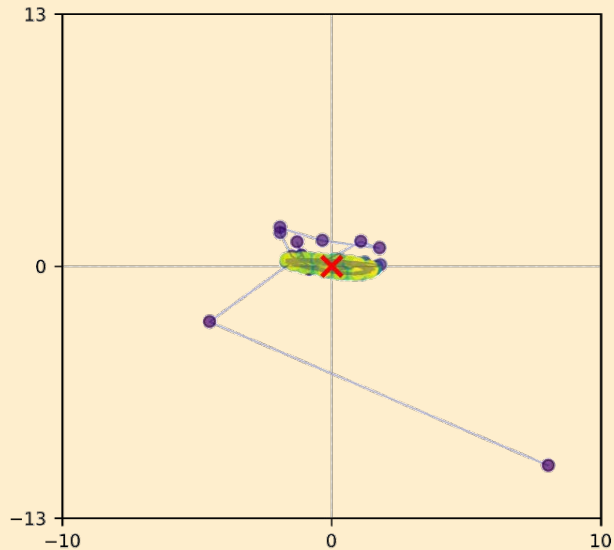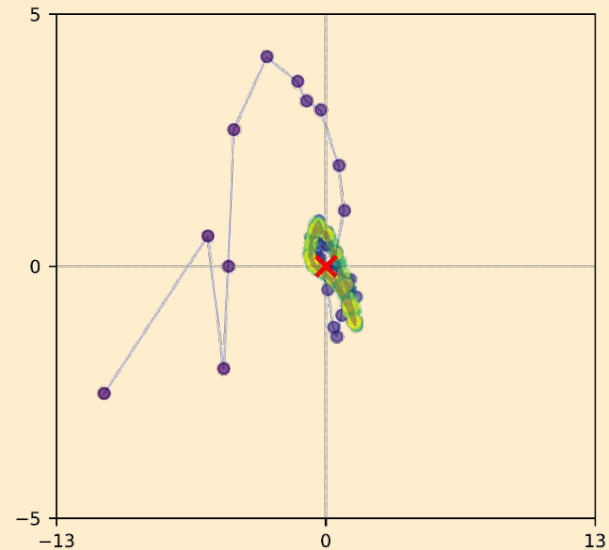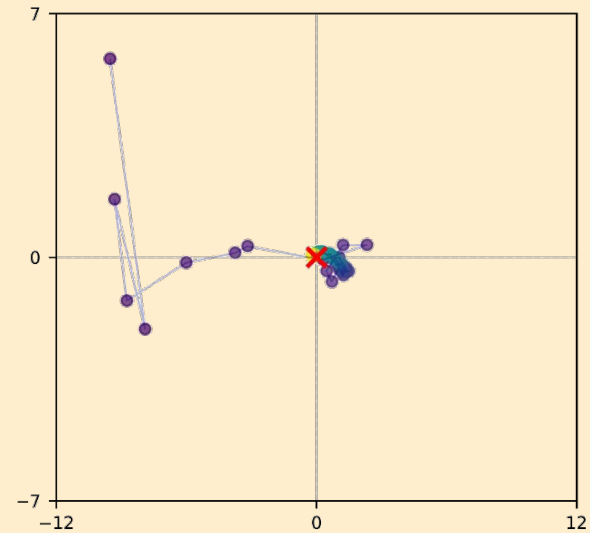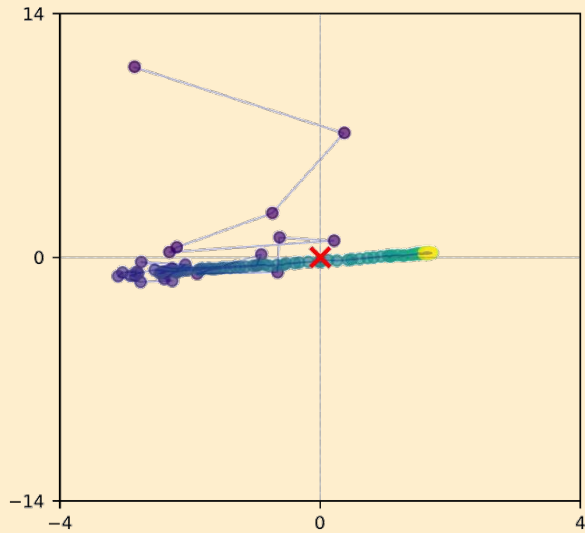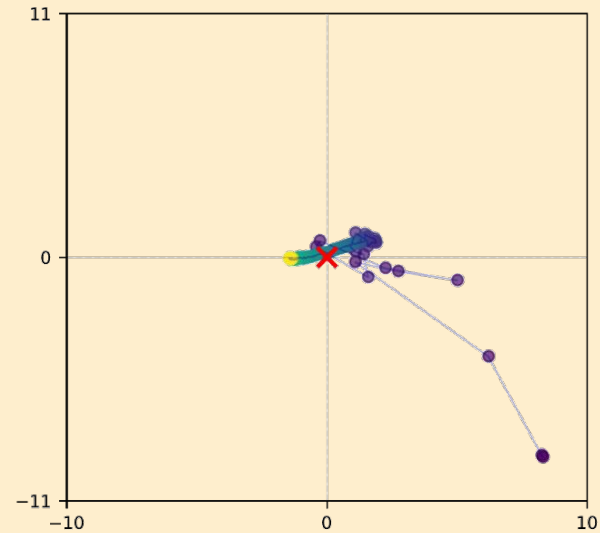
# Emergent Terminal Behaviors



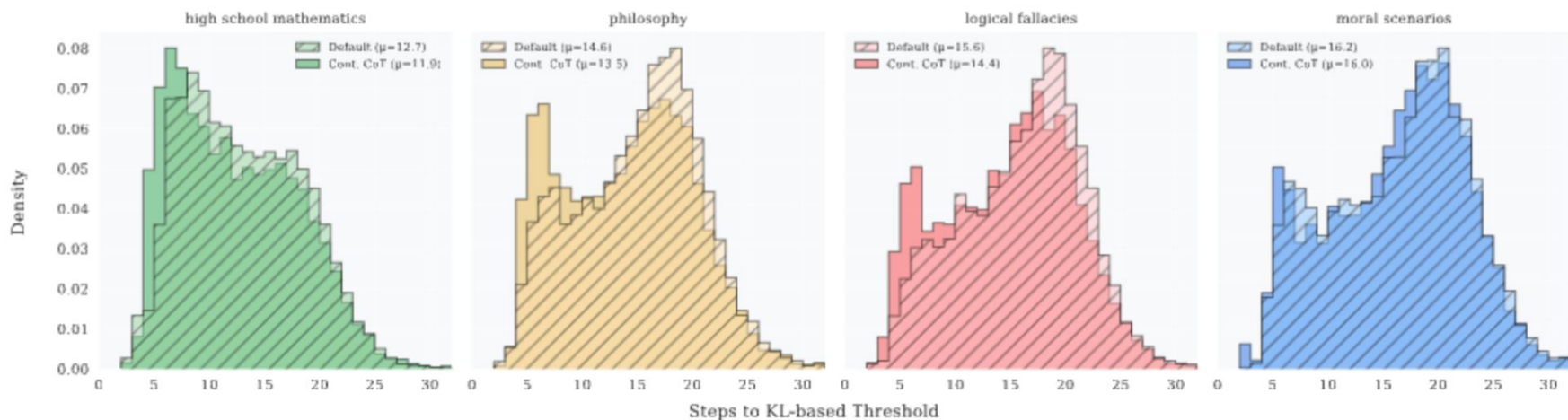Token: " wrong"
PC1-PC2

PC3-PC4

PC5-PC6

# Takeaways from Trajectories

- **Complexity emerges from pretraining**
- **Different terminal behaviors emerge from simple training objectives**
- **Harder to analyze model behavior -> requires representation analysis**

# Other Advantages of Recurrence

# The model is able to exit *per-token* without training

# Recurrent Depth actually simplifies LLMs

- Simplified speculative decoding
- Zero-shot per-token adaptive computation
- Simplified KV-cache sharing
- Simplified continuous chain of thought

Conclusions, Takeaways, the Future

# Conclusions, Takeaways, the Future

- Different paradigm to pre-train models that scales surprisingly far
- How do we get arbitrary extrapolation in compute?
- How to post-train?
- Is this a complementary path to scaling model performance? What is an apples-to-apples comparison to CoT?

# Questions