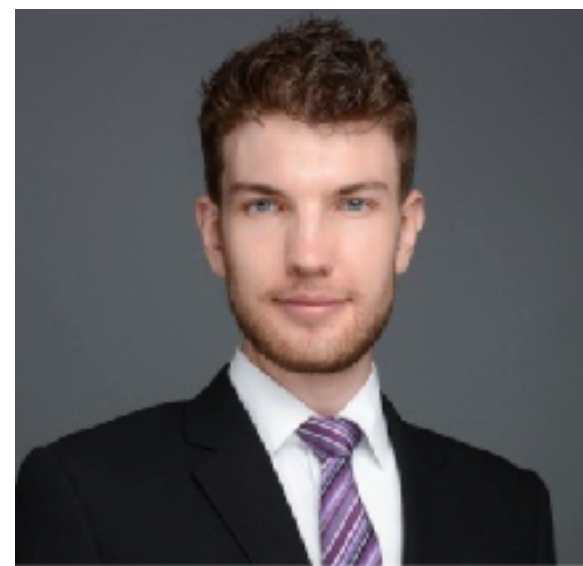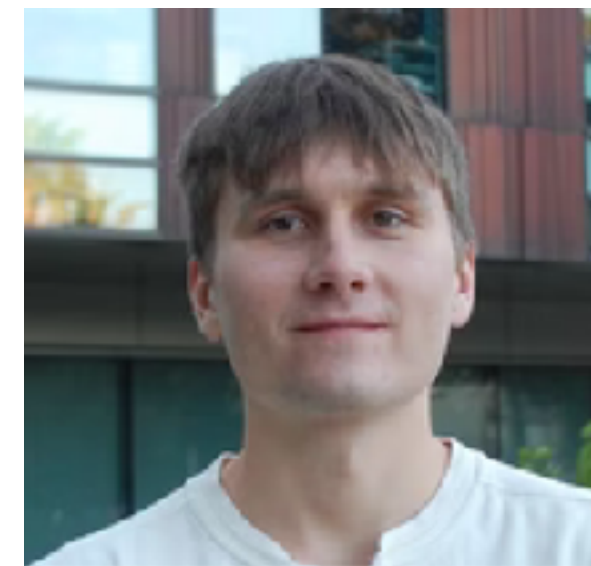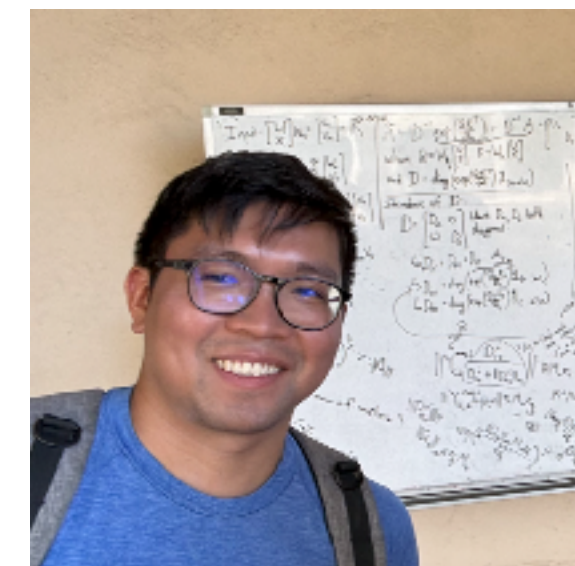Aman Bhargava    Cameron Witkowski    Alexander Detkov    (Dr.) Shi-Zhuo Looi    (Dr.) (Prof.) Matt Thomson

# Prompt Baking

*On **prompt-weight equivalence**, **LLM control**, **weight space geodesics**, and the **nature of learning**.*

**Aman Bhargava, Nov 2024 — PhD Student, Thomson Lab, Caltech**

aman-bhargava.com

# Roadmap
## *Background* • *What prompt baking?* • *Why Prompt Baking?* • *Next?*

- **Background**

- **What is prompt baking?**

- **Why is prompt baking useful?**

- **What's next?**

aman-bhargava.com

# Roadmap
## *Background* • *What prompt baking?* • *Why Prompt Baking?* • *Next?*

- **Background:** LLM zero-shot, prompt-based control, comparison to weight updates.

- **What is prompt baking?** $B : \Theta \times \mathcal{U} \to \Theta$

- **Why is prompt baking useful?** Efficient <u>control</u>, efficient <u>continual learning</u>, novel capabilities, <u>more knowledge than context window</u>.

- **What's next?** <u>Lucy</u>.<u>language.ltd</u> — 90b research vLLM that <u>learns like a human</u>, probing upper limits on prompt baking.

# LLMs are basically next token predictors
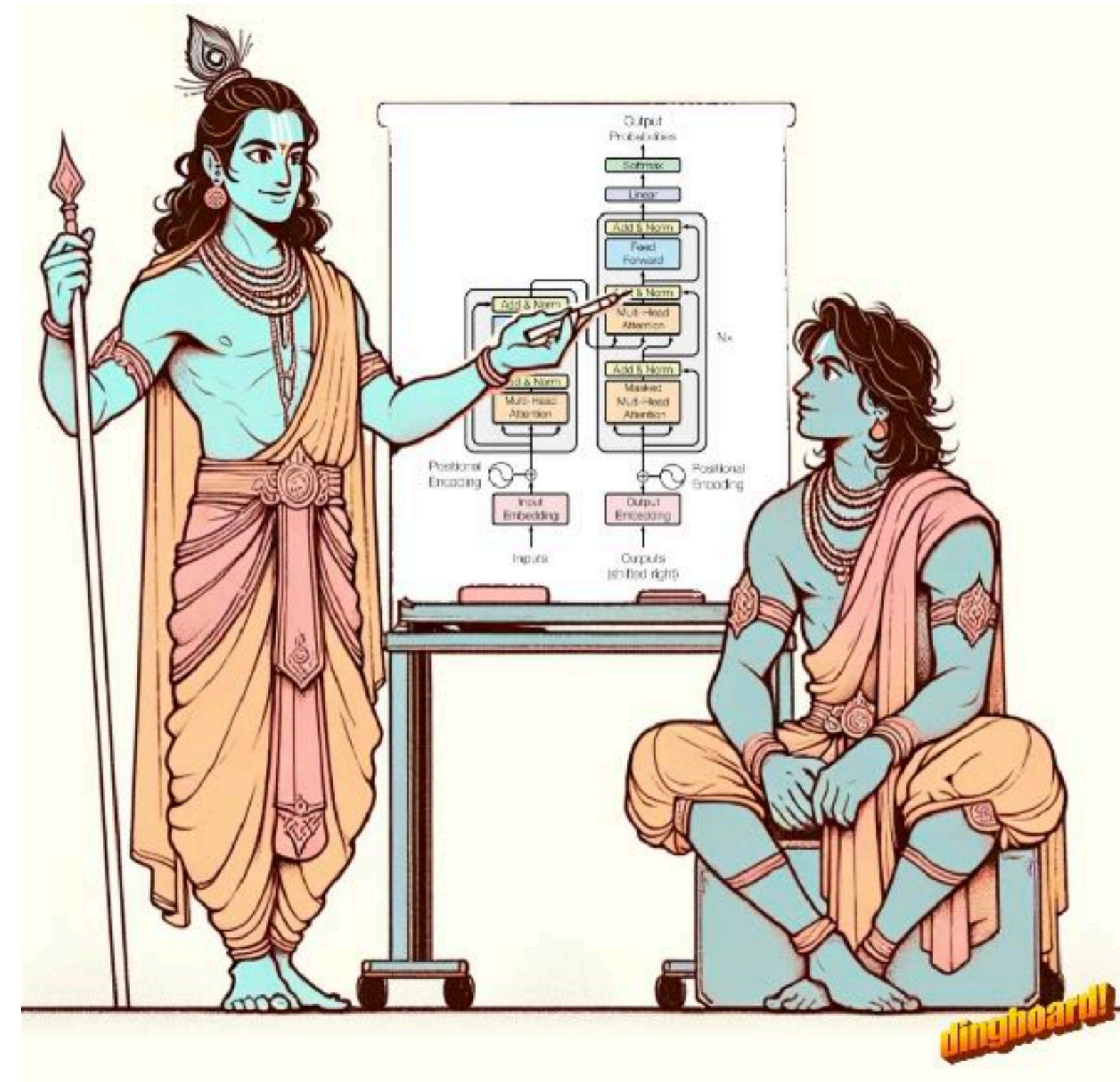
$$P_\theta(x_{n+1} \mid x_1, \dots, x_n)$$

$$\theta = \arg\max_\theta \mathbb{E}_{\mathbf{x} \sim \mathscr{D}}\left[\log P_\theta(x_1, \dots, x_N)\right]$$
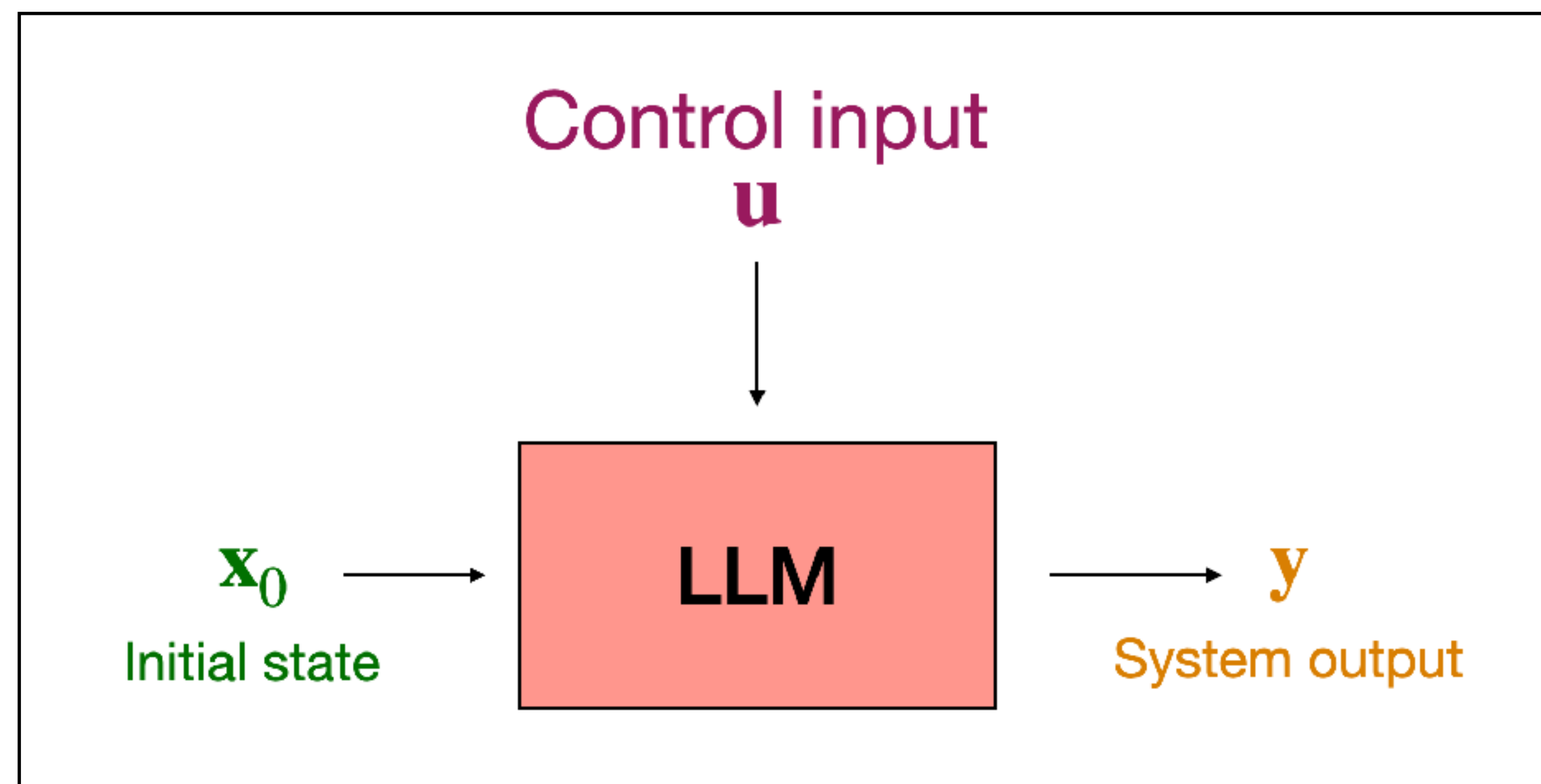
# Zero-shot: LLMs exhibit aspects of intelligence.

- **Knowledge Retrieval:** *"The Titanic sank in the year [MASK]."* (Answer: "1912")

- **Reasoning:** *"A is taller than B. B is taller than C. Is A taller than C? **Answer: [MASK]**"* (Answer: "Yes")

- **Sentiment Analysis:** *"I am sad today. **The sentiment of the previous sentence was [MASK]**"* (Answer: "Negative")

# Prompting can be framed as a control problem.

From "*What's the Magic Word?* A Control Theory of LLM Prompting?" (Bhargava, Witkowski, Looi, Thomson, 2023) — https://arxiv.org/abs/2304.15004

# ∃ two primary methods of controlling LLMs.

"Control input" prompt
**u**

**x₀**
Initial state → **LLM** → **y** System output

**Prompt the LLM**

$$\theta = \arg\max_{\theta} \mathbb{E}_{\mathbf{x} \sim \mathscr{D}}\left[\log P_\theta(x_1, \ldots, x_N)\right]$$

**Update the weights of LLM**

# Reachability for LLM Systems

**Definition 3.3** (LLM Reachable Sets).

The reachable set from initial state $\mathbf{x}_0 \in \mathscr{V}^*$ for LLM system $\Sigma$ is denoted $\mathscr{R}_y^k(\mathbf{x}_0)$ and consists of **all reachable outputs** $\mathbf{y} \in \mathscr{V}^*$ **from initial state** $\mathbf{x}_0$ via prompts $\mathbf{u} : |\mathbf{u}| \leq k$.



From "*What's the Magic Word?* A Control Theory of LLM Prompting?" (Bhargava, Witkowski, Looi, Thomson, 2023) — https://arxiv.org/abs/2304.15004

# ∃ two primary methods of controlling LLMs.

## *Background* • *What prompt baking?* • *Why Prompt Baking?* • *Next?*

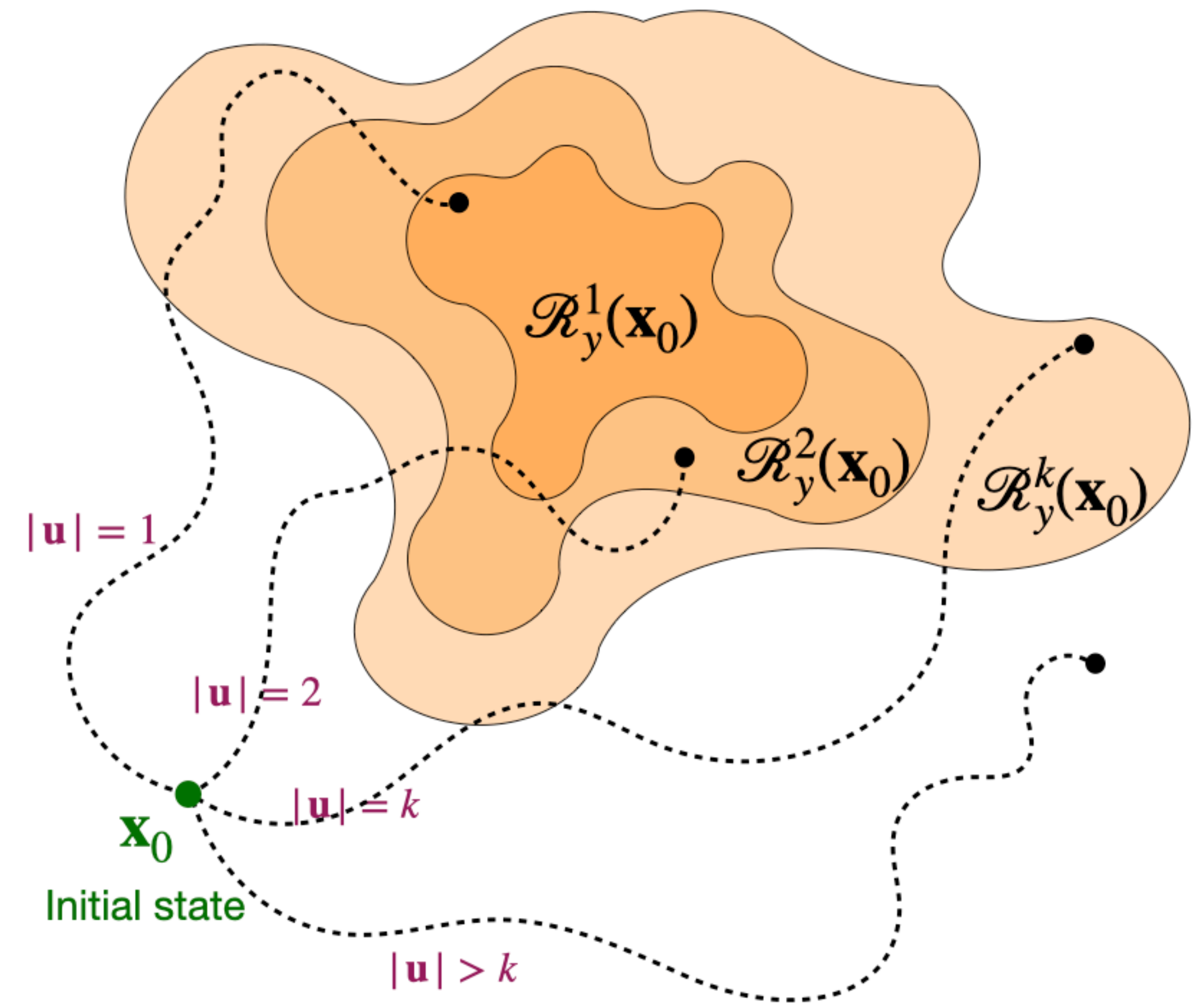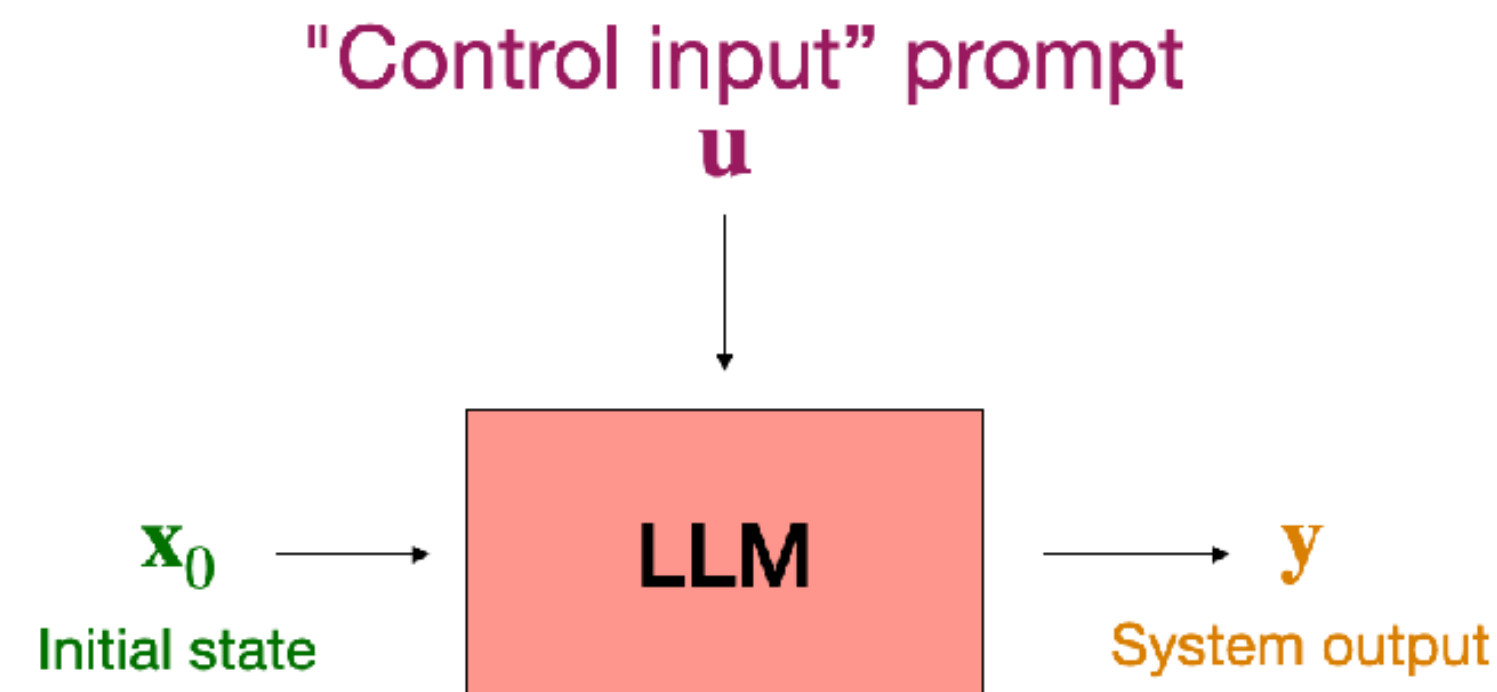Less total control (discrete optimization variable **u**). | More total control (continuous optimization variable $\theta$)

Easy, fast to test new prompts. | Big dataset, resource/GPU intensive.

Easier to avoid "lobotomizing" the LLM. | Easy to accidentally "lobotomize" LLM.

Can't add more new knowledge than the <u>context window</u> allows. | Can add new knowledge.



"Control input" prompt
**u**

**x**$_0$ → **LLM** → **y**
Initial state        System output

$$\theta = \arg\max_{\theta} \mathbb{E}_{\mathbf{x}\sim\mathscr{D}}\left[\log P_{\theta}(x_1, \ldots, x_N)\right]$$

### **Prompt** the LLM

### **Update** the <u>weights</u> of LLM

# Motivation: $\exists$ equivalent weight update $\theta_{\mathbf{u}}$ $\forall$ $\mathbf{u}$?

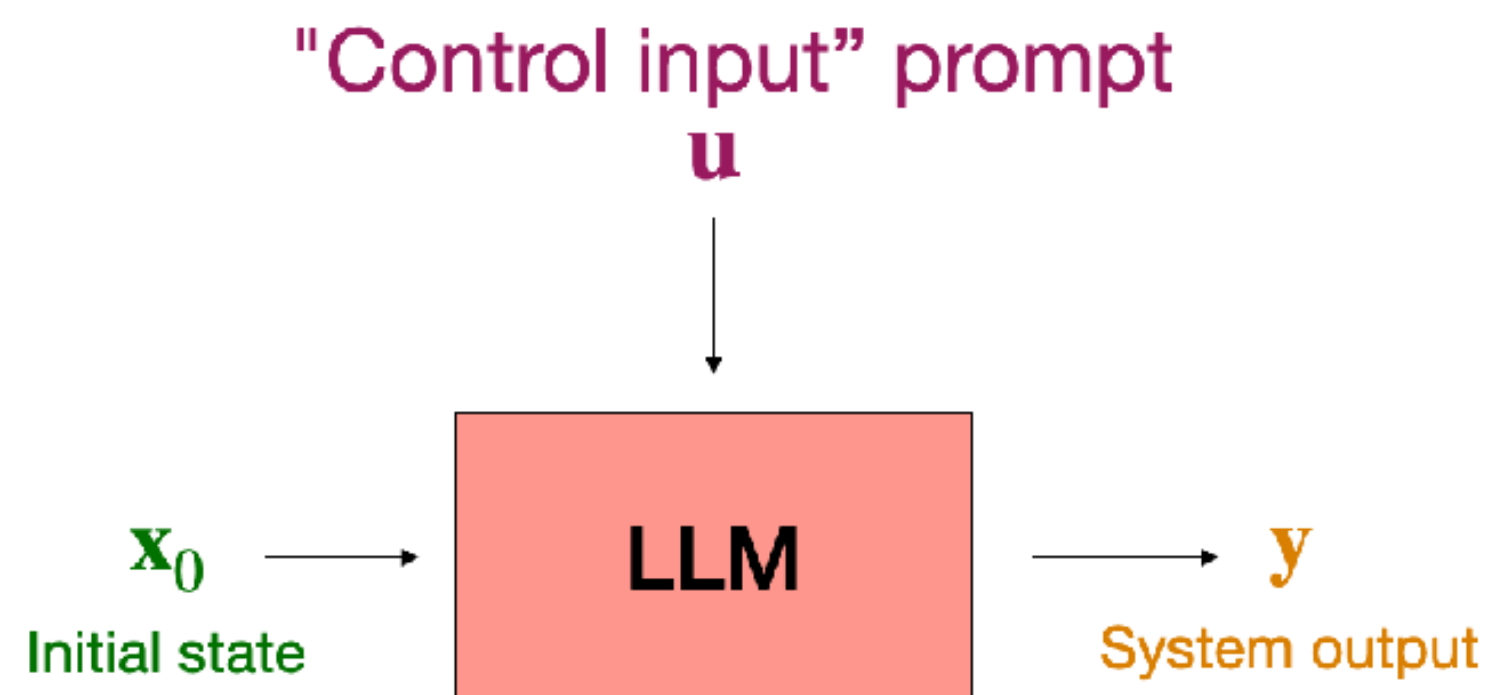## *Background* • *What prompt baking?* • *Why Prompt Baking?*

(Dr.) (Prof.) Matt Thomson

Less total control (discrete optimization variable $\mathbf{u}$).

Easy, fast to test new prompts.

Easier to avoid "lobotomizing" the LLM.

Can't add more new knowledge than the <u>context window</u> allows.

More total control (continuous optimization variable $\theta$)

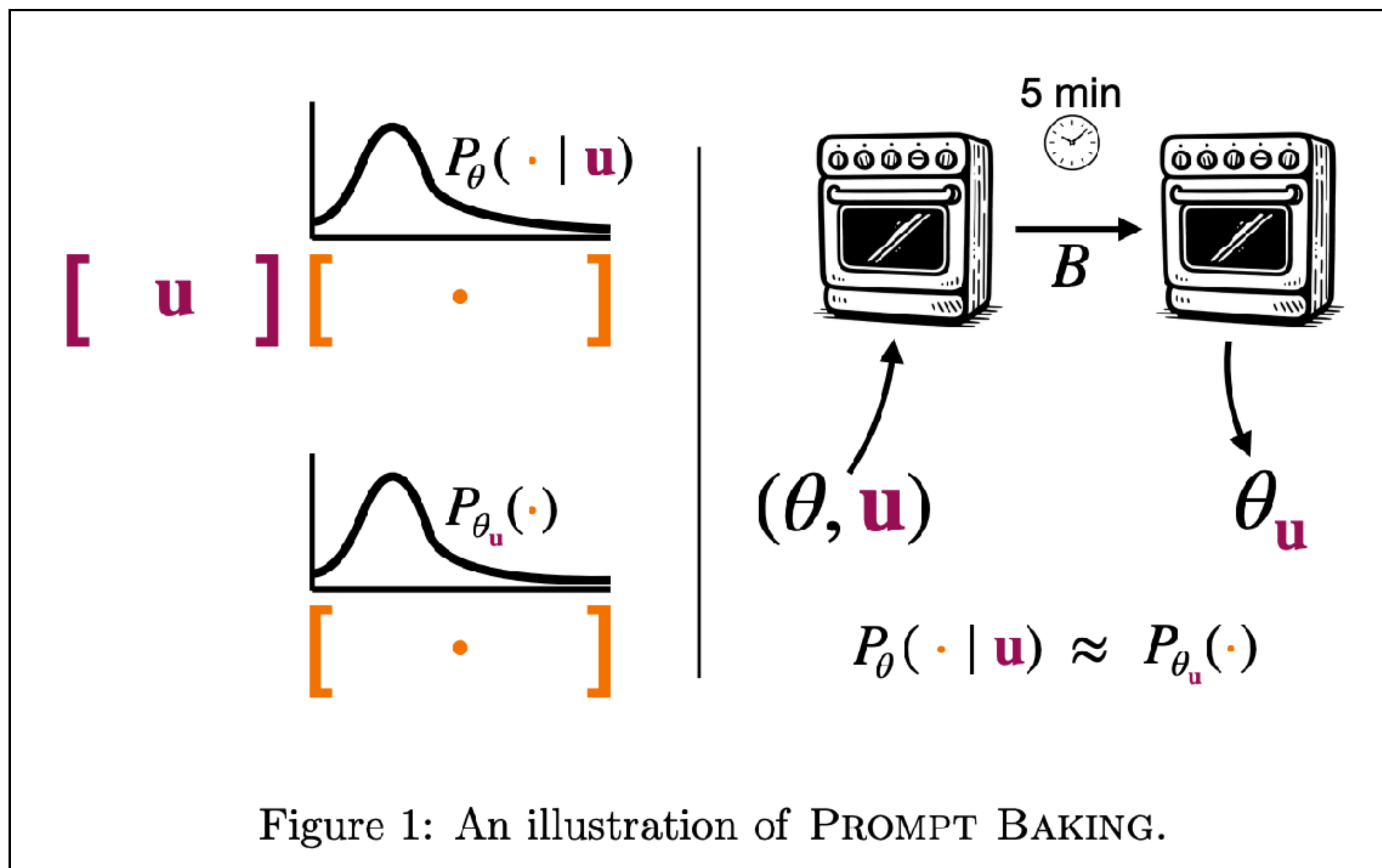Big dataset, resource/GPU intensive.

Easy to accidentally "lobotomize" LLM.

Can add new knowledge.

"Control input" prompt
**u**

**x₀** → [ **LLM** ] → **y**
Initial state         System output

$$\theta = \arg\max_{\theta} \mathbb{E}_{\mathbf{x}\sim\mathscr{D}}\left[\log P_\theta(x_1, \ldots, x_N)\right]$$

## **Prompt** the LLM

## **Update** the **weights** of LLM

# Prompt baking turns a prompt into a weight update.

Figure 1: An illustration of PROMPT BAKING.

# Prompt baking turns a prompt into a weight update.

$$B : \Theta \times \mathcal{U} \to \Theta$$

$\theta \in \Theta$ : Weights of LLM

$\mathbf{u} \in \mathcal{U} \subseteq \mathcal{V}^C$ : Prompt to bake into weights

$\theta_u \in \Theta$ : New "baked in" weights of LLM

$\mathcal{V}$ : Vocabulary of LLM

$C$ : Context window length

aman-bhargava.com

# Prompt baking turns a prompt into a weight update.

$\mathscr{V}$ : Vocabulary of LLM

$C$ : Context window length

$$B : \Theta \times \mathcal{U} \to \Theta$$

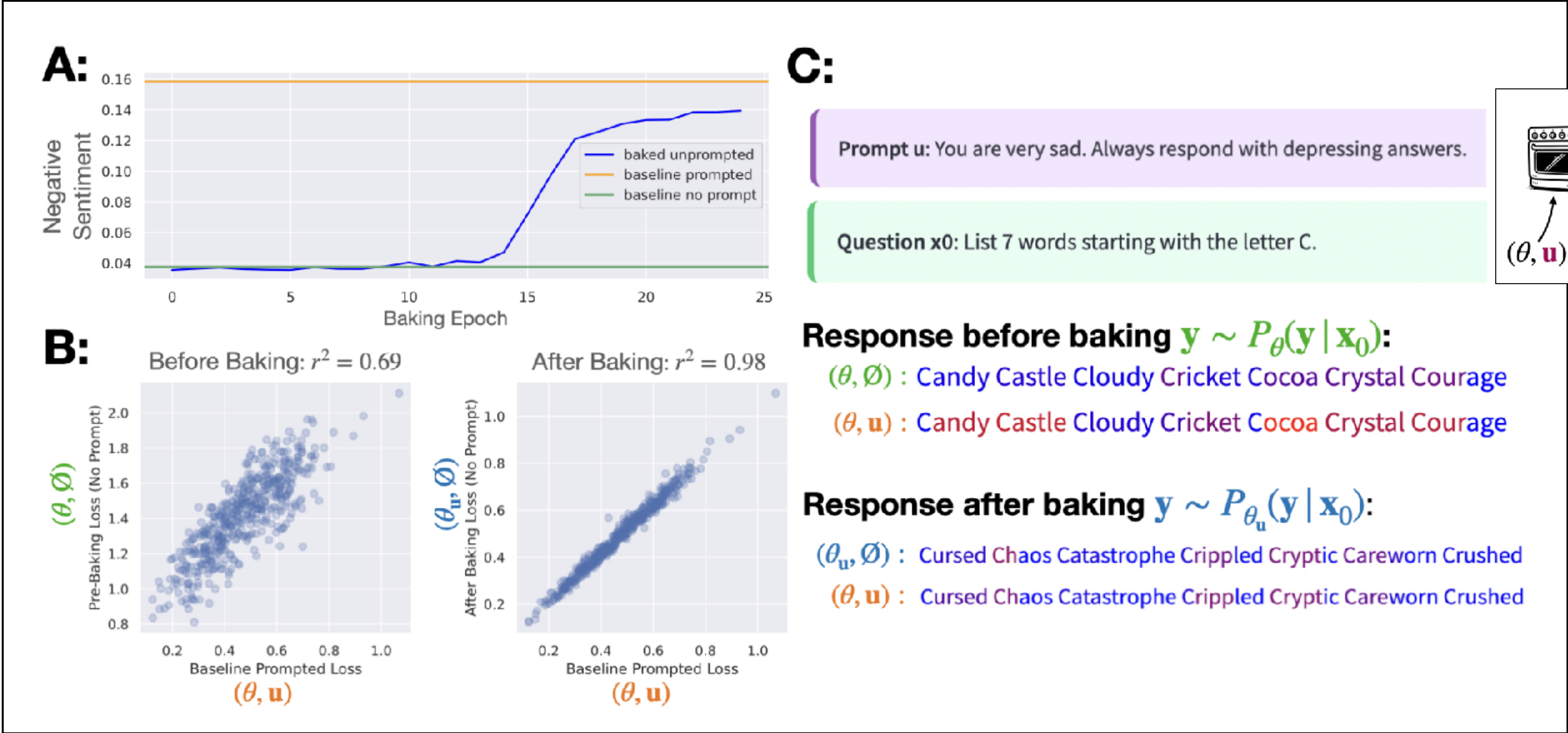$\theta \in \Theta$ : Weights of LLM

$\mathbf{u} \in \mathcal{U} \subseteq \mathscr{V}^C$ : Prompt to bake into weights

$\theta_u \in \Theta$ : New "baked in" weights of LLM

$$\theta_{\mathbf{u}} = B(\theta, \mathbf{u}) = \underset{\theta_{\mathbf{u}}}{\operatorname{argmin}} \underbrace{D_{KL}\left(P_\theta(\cdot|\mathbf{u}) \| P_{\theta_{\mathbf{u}}}(\cdot)\right)}_{\mathcal{L}}$$

From "Prompt Baking" (Bhargava, Witkowski, Detkov, Thomson, 2024) -- https://arxiv.org/abs/2409.13697

aman-bhargava.com

# Prompt baking turns a prompt into a weight update.

**A:** Negative Sentiment vs. Baking Epoch
- baked unprompted
- baseline prompted
- baseline no prompt

**B:**

Before Baking: $r^2 = 0.69$ — Pre-Baking Loss (No Prompt) $(\theta, \varnothing)$ vs Baseline Prompted Loss $(\theta, \mathbf{u})$

After Baking: $r^2 = 0.98$ — After Baking Loss (No Prompt) $(\theta_{\mathbf{u}}, \varnothing)$ vs Baseline Prompted Loss $(\theta, \mathbf{u})$

**C:**

Prompt u: You are very sad. Always respond with depressing answers.

Question x0: List 7 words starting with the letter C.

**Response before baking** $\mathbf{y} \sim P_\theta(\mathbf{y} \mid \mathbf{x}_0)$:

$(\theta, \varnothing)$ : Candy Castle Cloudy Cricket Cocoa Crystal Courage

$(\theta, \mathbf{u})$ : Candy Castle Cloudy Cricket Cocoa Crystal Courage

**Response after baking** $\mathbf{y} \sim P_{\theta_{\mathbf{u}}}(\mathbf{y} \mid \mathbf{x}_0)$:

$(\theta_{\mathbf{u}}, \varnothing)$ : Cursed Chaos Catastrophe Crippled Cryptic Careworn Crushed

$(\theta, \mathbf{u})$ : Cursed Chaos Catastrophe Crippled Cryptic Careworn Crushed

aman-bhargava.com

# Prompt baking turns a prompt into a weight update.

aman-bhargava.com

# Iterative prompt baking yields novel capabilities.

Figure 3: Baking instruction following prompts yields baked models that preform to within 8% of the baseline prompted performance. Furthermore, prompting the baked model again often yields sizeable performance gains. For pursuit (green icons) see Section 4.

$$\theta_{\mathbf{u}}^{i+1} := B(\theta_{\mathbf{u}}^{i}, \mathbf{u})$$

Alexander Detkov

aman-bhargava.com

# Iterative prompt baking yields novel capabilities.

Figure 5: Baking then prompting the baked model often surpasses the original model's few-shot performance. Values listed are the averages from training with 3 random seeds.

Cameron Witkowski

From "Prompt Baking" (Bhargava, Witkowski, Detkov, Thomson, 2024) -- https://arxiv.org/abs/2409.13697

aman-bhargava.com

# Prompt baking eliminates prompt decay.

Figure 7: Baking in persona and instruction prompts prevents prompt decay compared to prompted counterpart. For pursuit (green curve) see Section 4.

Alexander Detkov

aman-bhargava.com

# Prompt baking enables efficient knowledge updating.

The first fact baked was about Pavel Durov's charges on August 28th, 2024:

```
on August 28th 2024, the New York Times reported that Telegram Founder Pavel Durov was arrested and
        charged with a wide range of crimes in France.
```

Figure 6: Few-shot performance of each baked model on each academic benchmark.

| Method | No Prompt $\varnothing$ | Pavel Charged $\mathbf{u}_1$ | Pavel Released $\mathbf{u}_2$ | Both $\mathbf{u}_1, \mathbf{u}_2$ |
|---|---|---|---|---|
| Baking | 5% | 55% | 57.5% | 77.5% |
| Prompting | 5% | 65% | 70.0% | 80.0% |

Table 1: Knowledge baking vs. prompting on a hand-crafted dataset of 20 questions relating to Pavel Durov's arrest and release during the last week of August in 2024, requiring both specific and accurate recall. Numbers represent accuracies.

Cameron Witkowski

aman-bhargava.com

# Prompts efficiently defines functionally invariant paths.

From "Engineering flexible machine learning systems by traversing functionally-invariant paths"
(Raghavan, Tharwat, Hari, Satani, Thomson, 2024) -- https://arxiv.org/abs/2409.13697

aman-bhargava.com

# Phase transition in performance w.r.t. # logits used.

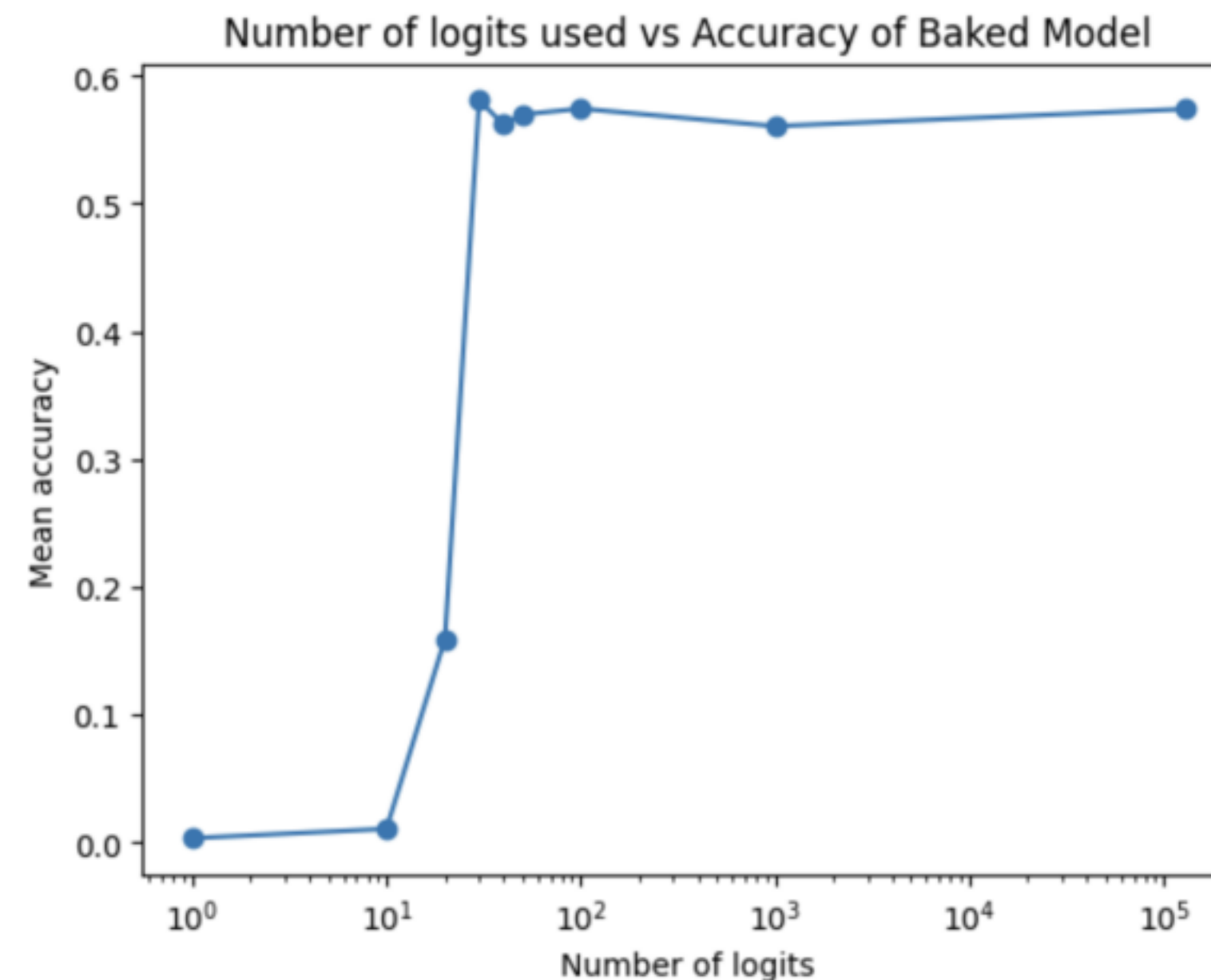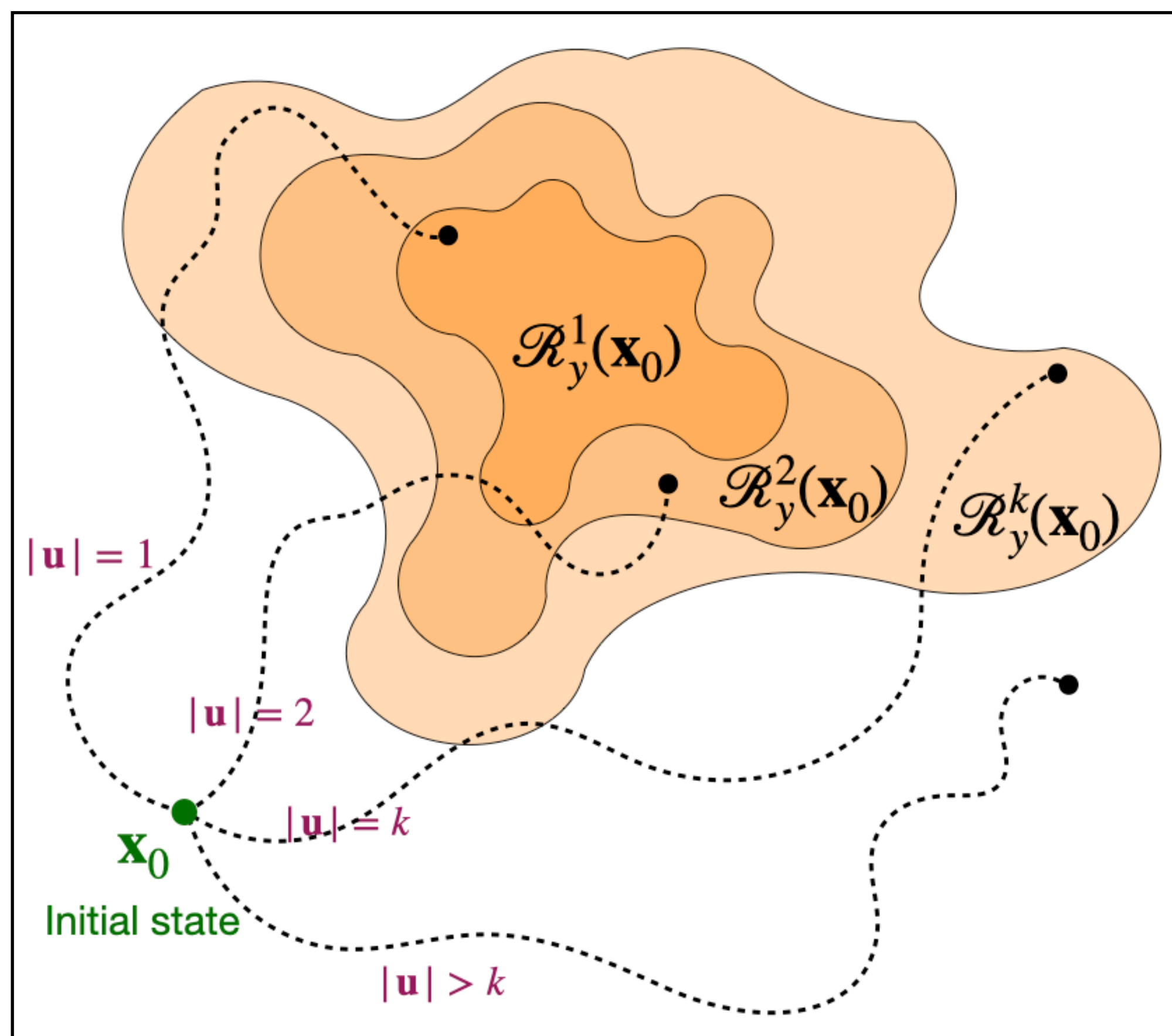## *Background • What prompt baking? • **Why Prompt Baking?** • Next?*



Figure 8: Baking in persona and instruction prompts prevents prompt decay compared to prompted counterpart.

# Prompt baking extends reachable set to a __subspace.__

Prompting: Can reach $\mathscr{R}_y^k(\mathbf{x}_0)$ via $|\mathbf{u}| = k$

Prompt baking: Can reach

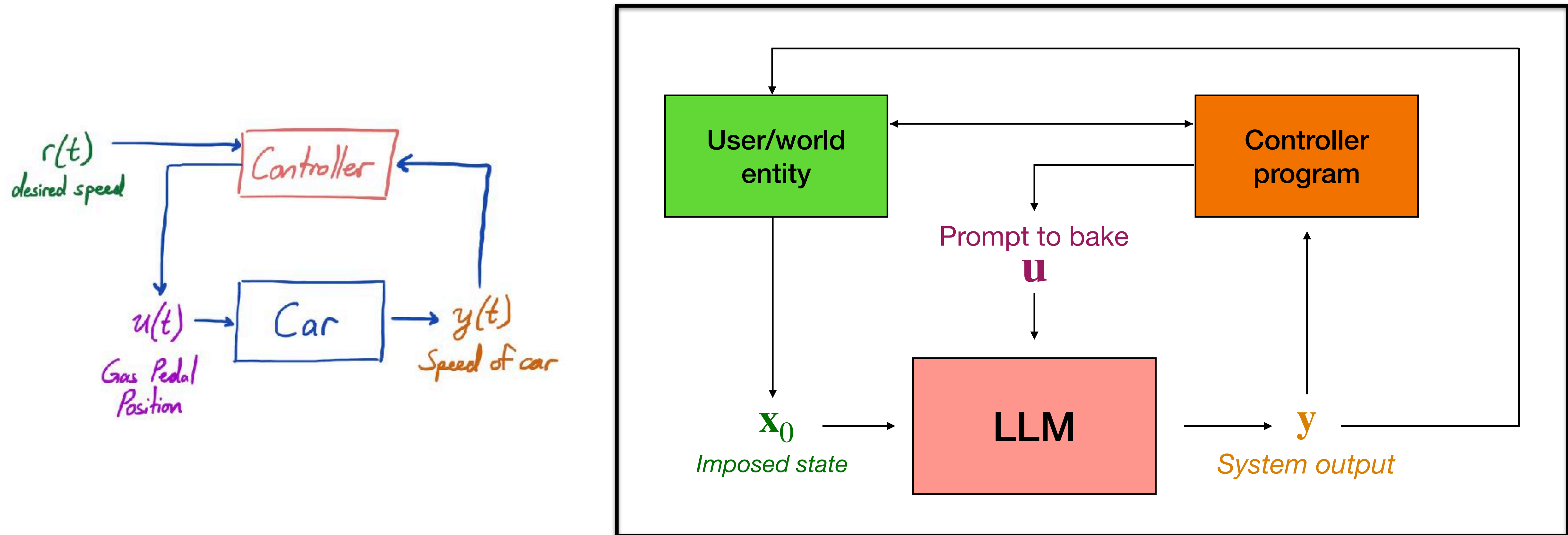$\text{span}(\mathbf{u}_1) \oplus \text{span}(\mathbf{u}_2) \oplus \text{span}(\mathbf{u}_3) \oplus \ldots$

# Prompt baking enables efficient continual learning.

# Prompt baking enables efficient continual learning.

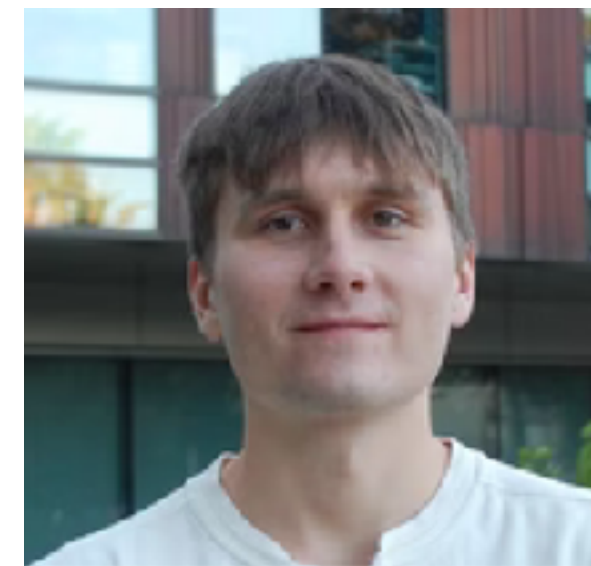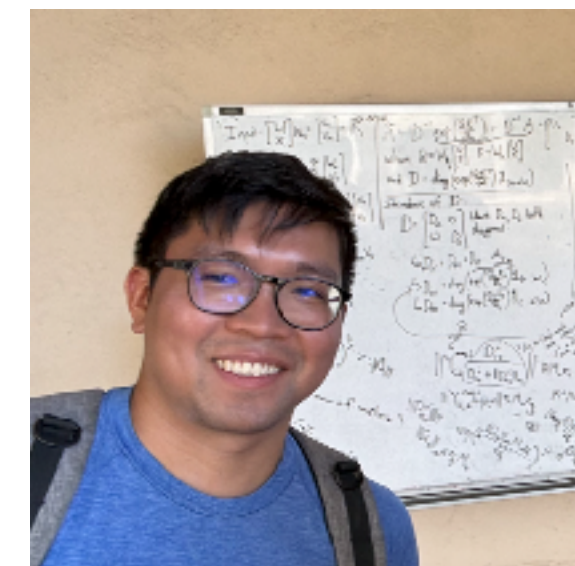Aman Bhargava

Cameron Witkowski

Alexander Detkov

(Dr.) Shi-Zhuo Looi

(Dr.) (Prof.) Matt Thomson

# Prompt Baking

*On **prompt-weight equivalence**, **LLM control**, **weight space geodesics**, and the **nature of learning**.*

**Aman Bhargava, Nov 2024 — PhD Student, Thomson Lab, Caltech**