



# Predictive, Scalable and Interpretable knowledge tracing on structured domains



Hanqi Zhou, Robert Bamler, Charley M. Wu\*, Álvaro Tejero-Cantero\* (\*equal contribution)

[hanqi.zhou@uni-tuebingen.de](mailto:hanqi.zhou@uni-tuebingen.de)

Department of Computer Science,

Cluster of Excellence "Machine Learning for Science", Tübingen AI Center,

University of Tübingen

# Tracing knowledge over time

We aim to improve self-directed learning

- by estimating learner knowledge, and providing the right learning materials

Type the answer



32

CHECK

1

2

3

4

5

6

7

8

9

0

✖

Select the answer

$$34 - 0 = \square$$

7

31

34

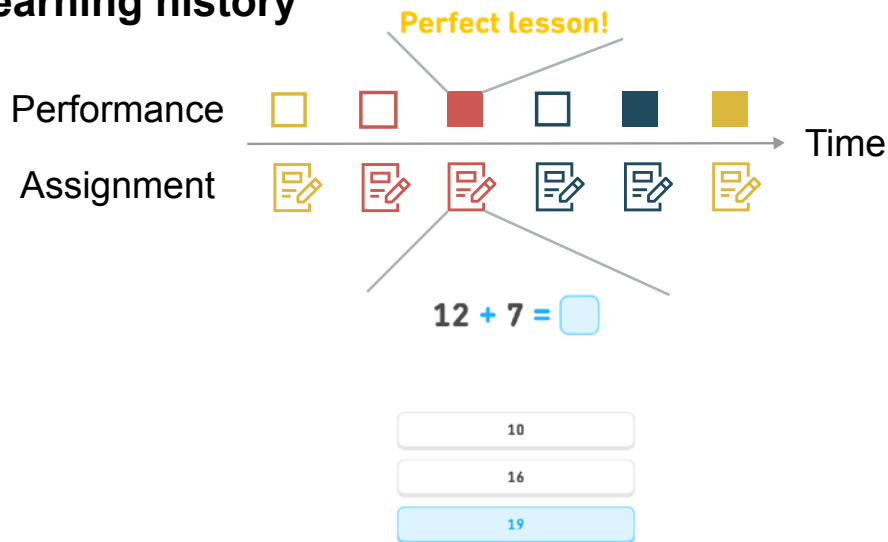
CONTINUE

source: Duolingo (math learning)

- We observe learning histories of each learner

- We observe learning histories of each learner

## Learning history



source: Duolingo (math learning)

- We observe learning histories of each learner

## Learning history



## Knowledge structure

Knowledge concept (KC)



source: Duolingo (math learning)

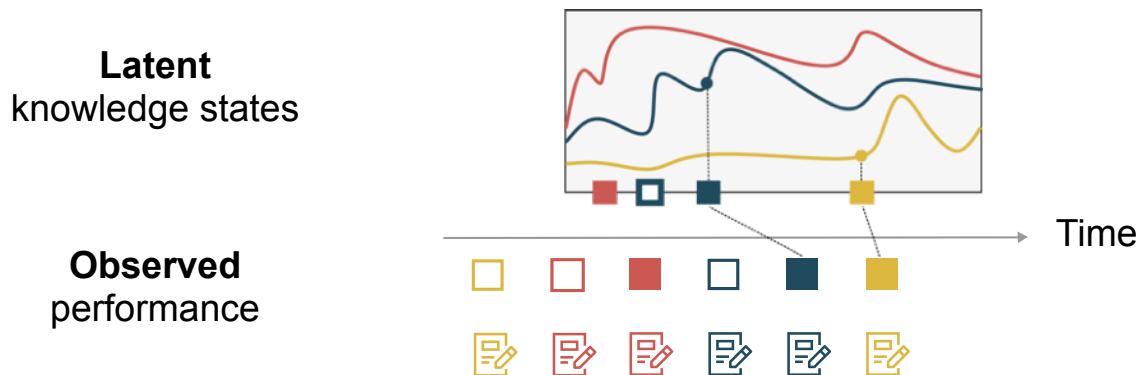
- We observe learning histories of each learner
- We aim to estimate a learner's knowledge states and predict future performance



- We observe learning histories of each learner
- We aim to estimate a learner's knowledge states and predict future performance

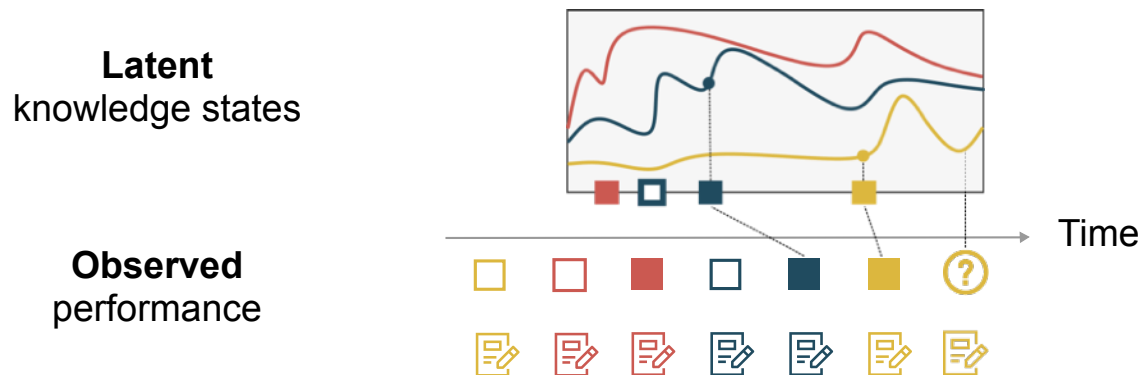


- We observe learning histories of each learner
- We aim to estimate a learner's knowledge states and predict future performance



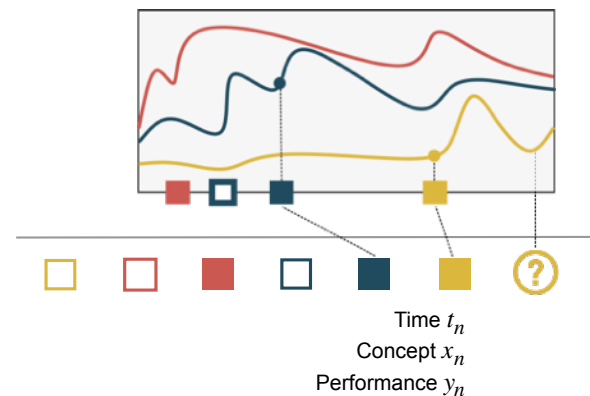


- We observe learning histories of each learner
- We aim to estimate a learner's knowledge states and predict future performance



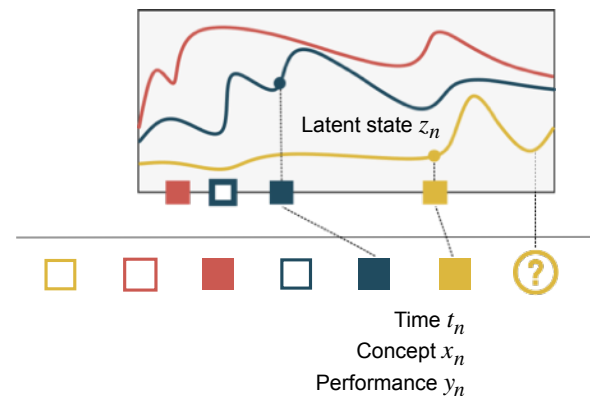
# Existing models just can't keep up!

- Learning history  $\mathcal{H}_{1:N}^{\ell} := \{t_n, x_n, y_n\}_{1:N}^l$  from learner  $\ell$ 
  - Interaction time  $t_n$
  - Knowledge concept index  $x_n \in \{1, \dots, K\}$
  - Learner's performance  $y_n \in [0, 1]$



# Existing models just can't keep up!

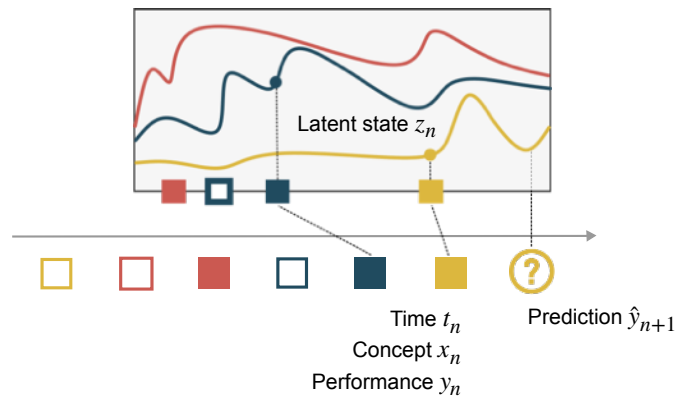
- Learning history  $\mathcal{H}_{1:N}^{\ell} := \{t_n, x_n, y_n\}_{1:N}^l$  from learner  $\ell$ 
  - Interaction time  $t_n$
  - Knowledge concept index  $x_n \in \{1, \dots, K\}$
  - Learner's performance  $y_n \in [0, 1]$
- Latent knowledge states  $\mathbf{z}_{1:N}^{\ell} = [z_1^{1:K}, \dots, z_n^{1:K}]^T$



# Existing models just can't keep up!

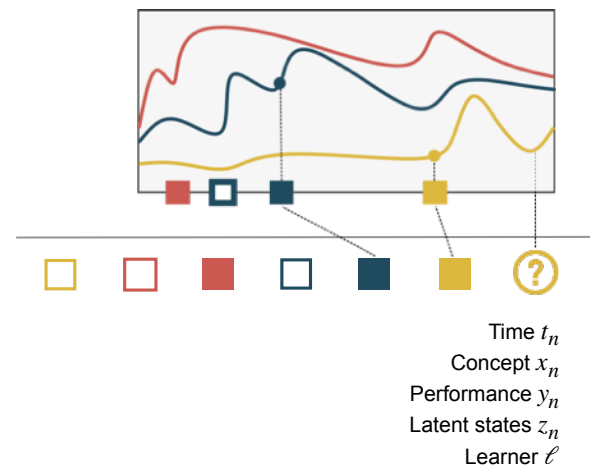


- Learning history  $\mathcal{H}_{1:N}^{\ell} := \{t_n, x_n, y_n\}_{1:N}^l$  from learner  $\ell$ 
  - Interaction time  $t_n$
  - Knowledge concept index  $x_n \in \{1, \dots, K\}$
  - Learner's performance  $y_n \in 0, 1$
- Latent knowledge states  $\mathbf{z}_{1:N}^{\ell} = [z_1^{1:K}, \dots, z_n^{1:K}]^T$
- Prediction  $\hat{y}_{n+1}$ 
  - $p(y_{n+1} = 1) = \text{sigmoid}(z_{n+1})$



Psychological methods: multi-factor regression

$$\begin{aligned} &f_{\theta}(z_{n+1} \mid \mathcal{H}_{1:n}, t_{n+1}, x_{n+1}) \\ &= \alpha \cdot \text{property}_{x_{n+1}} \\ &+ \beta \cdot \text{spacing}_{\mathcal{H}} \\ &+ \gamma \cdot \text{ability}_{\ell} \\ &+ \dots \end{aligned}$$

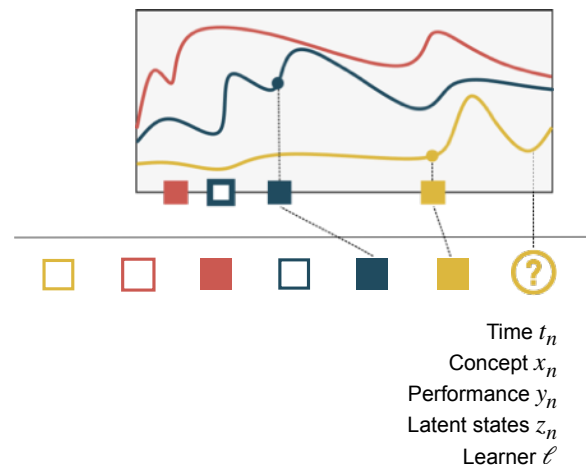


Psychological methods: multi-factor regression

$$\begin{aligned} & f_{\theta}(z_{n+1} \mid \mathcal{H}_{1:n}, t_{n+1}, x_{n+1}) \\ &= \alpha \cdot \boxed{\text{property}_{x_{n+1}}} \\ &+ \beta \cdot \text{spacing}_{\mathcal{H}} \\ &+ \gamma \cdot \text{ability}_{\ell} \\ &+ \dots \end{aligned}$$





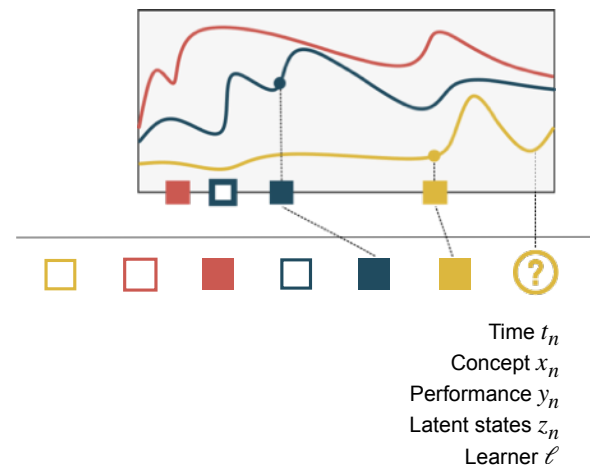
KC/assignment difficulty  
Correct/incorrect frequency



 Psychological methods: multi-factor regression

$$\begin{aligned} & f_{\theta}(z_{n+1} \mid \mathcal{H}_{1:n}, t_{n+1}, x_{n+1}) \\ &= \alpha \cdot \text{property}_{x_{n+1}} \\ &+ \beta \cdot \boxed{\text{spacing}_{\mathcal{H}}} \\ &+ \gamma \cdot \text{ability}_{\ell} \\ &+ \dots \end{aligned}$$

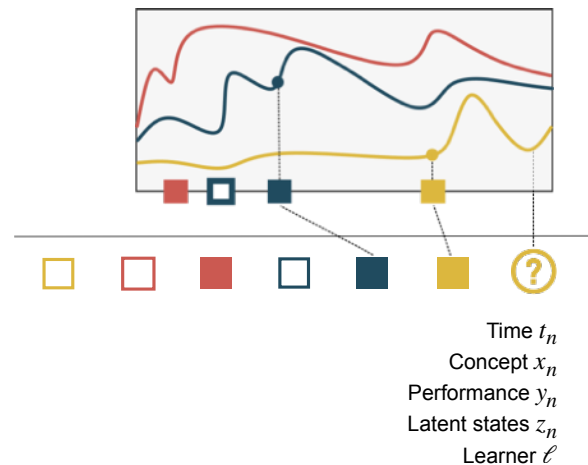
 Time duration  
 Recency effects





## Psychological methods: multi-factor regression

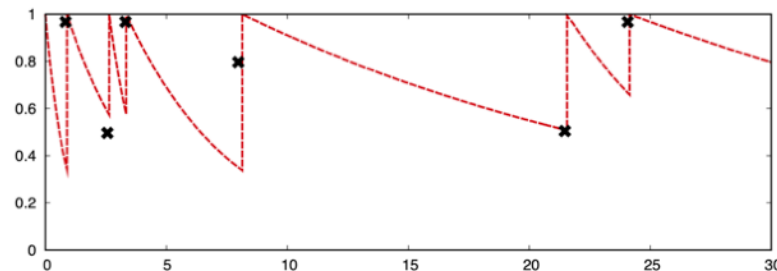
$$\begin{aligned} & f_{\theta}(z_{n+1} \mid \mathcal{H}_{1:n}, t_{n+1}, x_{n+1}) \\ &= \alpha \cdot \text{property}_{x_{n+1}} \\ &+ \beta \cdot \text{spacing}_{\mathcal{H}} \\ &+ \gamma \cdot \boxed{\text{ability}_{\ell}} \quad \text{Memory capacity} \\ &\quad \text{Guessing rate} \\ &+ \dots \end{aligned}$$





 Half-life regression (HLR; Settles & Seeder, 2016; Duolingo)

$$\begin{aligned} & f_{\theta}(z_{n+1} \mid \mathcal{H}_{1:n}, t_{n+1}, x_{n+1}) \\ &= \alpha \cdot \text{property}_{x_{n+1}} \quad \text{Correct/incorrect frequency} \\ &+ \beta \cdot \text{spacing}_{\mathcal{H}} \quad \text{Time duration} \end{aligned}$$



(b) 30-day student-word learning trace and predicted forgetting curve


Settles, B., & Meeder, B. (2016, August). A trainable spaced repetition model for language learning. In Proceedings of the 54th annual meeting of the association for computational linguistics (volume 1: long papers) (pp. 1848-1858).

# Existing models just can't keep up!



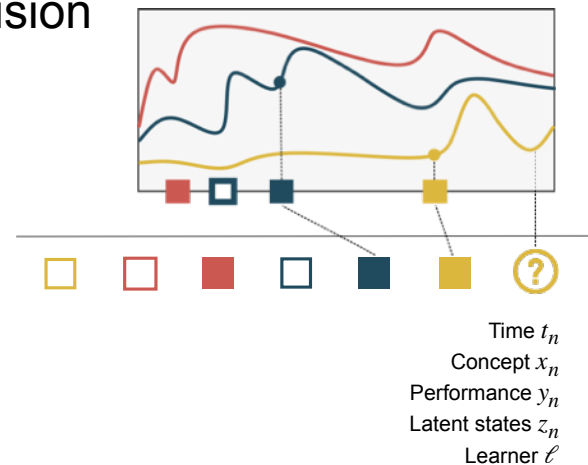
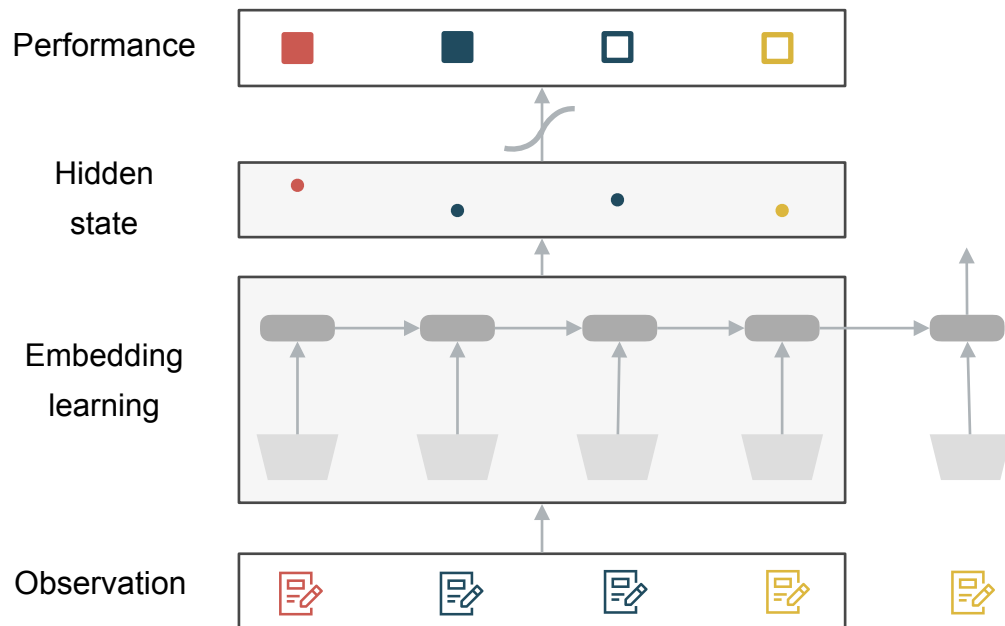
 Psychological methods: multi-factor regression

 Diagnosis: explicit modeling of performance factors

 Inflexibility: the amount of parameters increase as the learners/concepts increase

# Existing models just can't keep up!

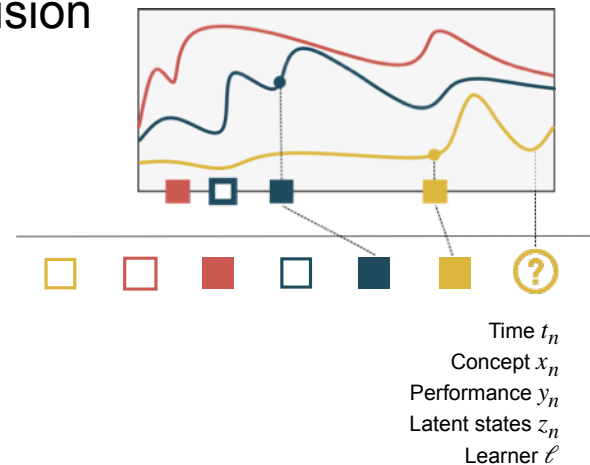
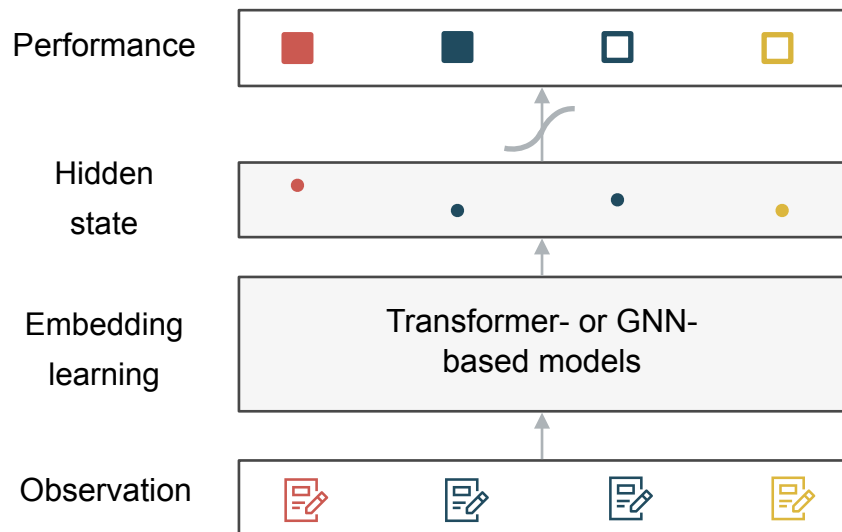
## Deep learning models: embedding learning and fusion



Piech, C., Bassen, J., Huang, J., Ganguli, S., Sahami, M., Guibas, L. J., & Sohl-Dickstein, J. (2015). Deep knowledge tracing. Advances in neural information processing systems, 28.

# Existing models just can't keep up!

## Deep learning models: embedding learning and fusion



Choi, Y., Lee, Y., Cho, J., Baek, J., Kim, B., Cha, Y., ... & Heo, J. (2020, August). Towards an appropriate query, key, and value computation for knowledge tracing. In *Proceedings of the seventh ACM conference on learning@ scale* (pp. 341-344).

Nakagawa, H., Iwasawa, Y., & Matsuo, Y. (2019, October). Graph-based knowledge tracing: modeling student proficiency using graph neural network. In *IEEE/WIC/ACM International Conference on Web Intelligence* (pp. 156-163).

# Existing models just can't keep up!



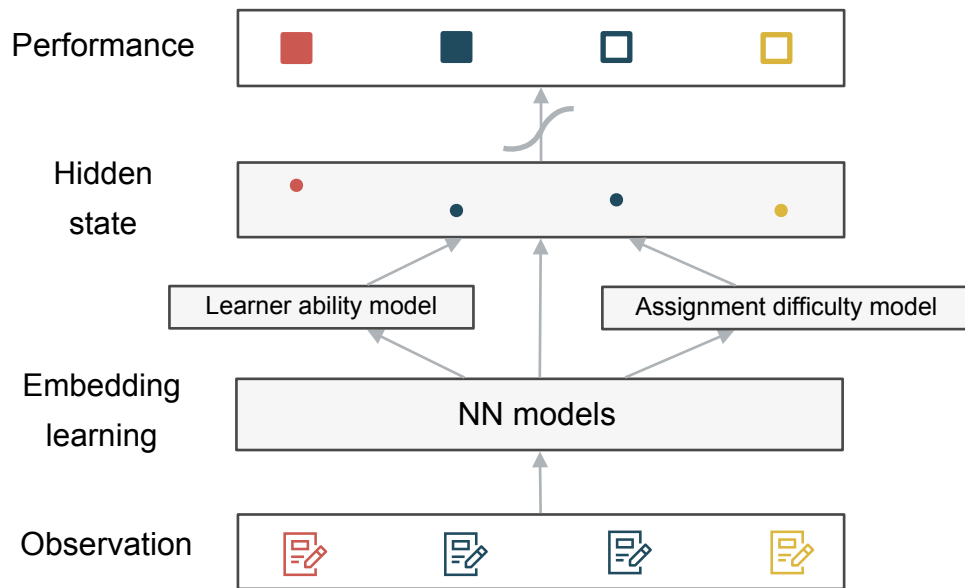
Deep learning models: embedding learning and fusion

High-capacity: can handle high-dimensional feature and large data

Interpretability: but what do these features mean?

# Existing models just can't keep up!

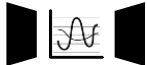
- A little more efforts: deep learning + psychology




$$\begin{aligned} p(y_{n+1}^{\text{correct}} \mid \mathcal{H}_{1:n}, t_{n+1}, x_{n+1}) \\ = \alpha \cdot \text{property}_{x_{n+1}} \\ + \beta \cdot \text{ability}_{\ell} \\ + \gamma \cdot \text{spacing}_{\mathcal{H}} \end{aligned}$$

**QIKT:** Chen, J., Liu, Z., Huang, S., Liu, Q., & Luo, W. (2023, June). Improving interpretability of deep sequential knowledge tracing models with question-centric cognitive representations. In *Proceedings of the AAAI Conference on Artificial Intelligence* (Vol. 37, No. 12, pp. 14196-14204).



- Computational modeling  :
- A flexible architecture with
- Identifiable interpretability





- Computational modeling  :
  - A flexible architecture with
  - Identifiable interpretability
- Real-world scenario:
  - Small data regime and real-time adaptation



# How do humans learn in structured domains?

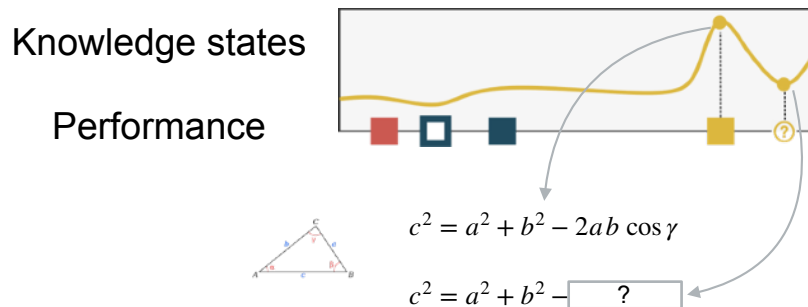
---



-  Human learners **forget** over time but reinforce knowledge through practice
-  Prerequisites **structure** the domain by relating KCs

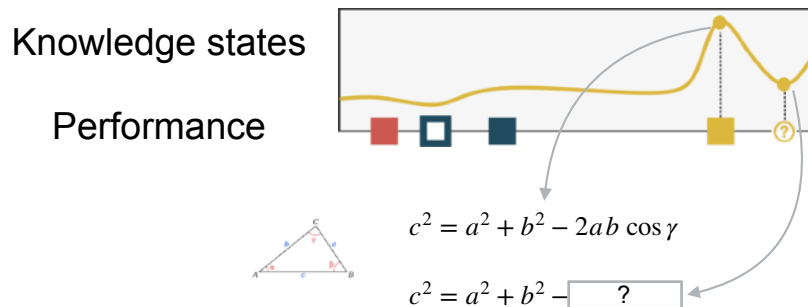
- 🧠 Human learners **forget** over time but reinforce knowledge through practice
- 🌐 Prerequisites **structure** the domain by relating KCs

## Learning process

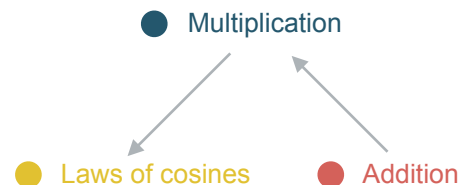


- 🧠 Human learners **forget** over time but reinforce knowledge through practice
- 🌐 Prerequisites **structure** the domain by relating concepts

## Learning process



## Prerequisites



# We jointly model learner traits and prerequisites



Observed  
performance  $y$



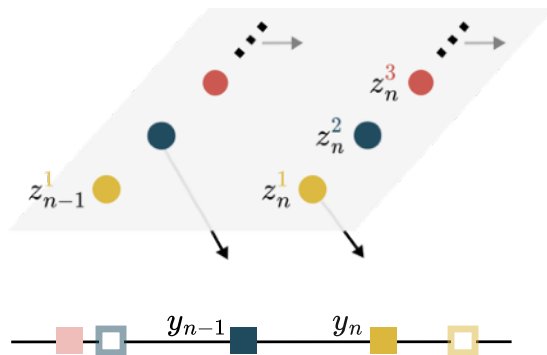
# We jointly model learner traits and prerequisites



Latent  
knowledge states  $z$



Observed  
performance  $y$



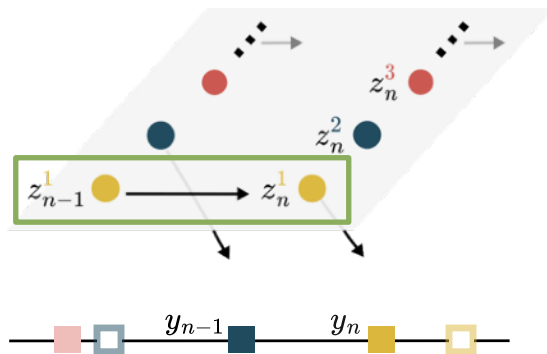
# We jointly model learner traits and prerequisites



Latent  
knowledge states  $z$



Observed  
performance  $y$



$$z_n^1 = f_{\mathcal{S}}^*(z_{n-1}^1; \theta_s)$$

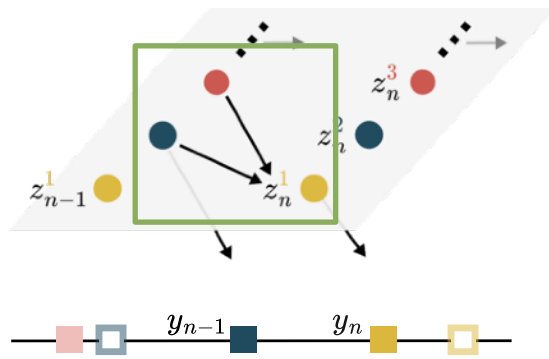
# We jointly model learner traits and prerequisites



Latent  
knowledge states  $z$



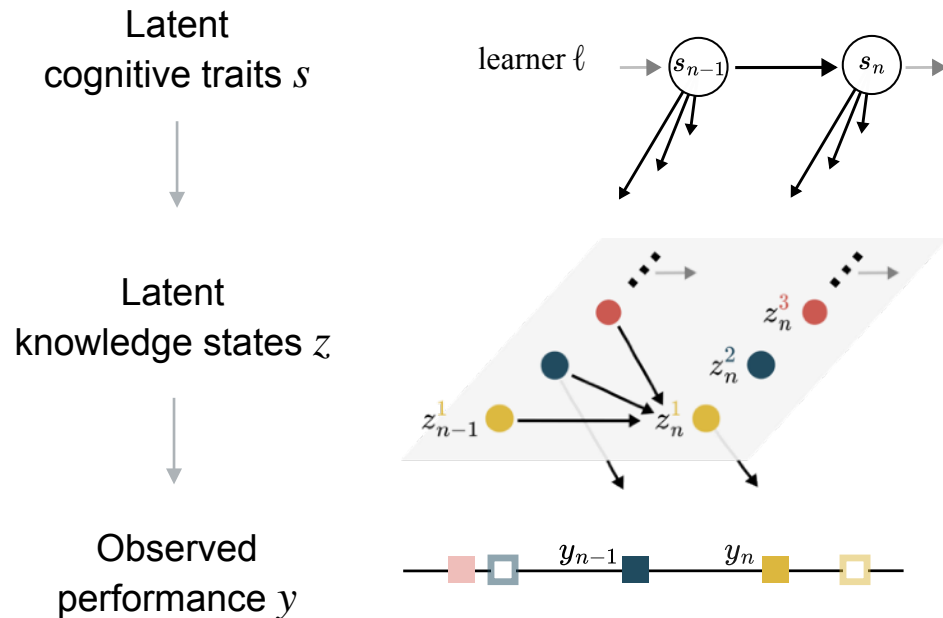
Observed  
performance  $y$



$$z_n^1 = f_{\text{graph}}(z_{n-1}^2, z_{n-1}^3; \theta_G)$$



# We jointly model learner traits and prerequisites



# We jointly model learner traits and prerequisites



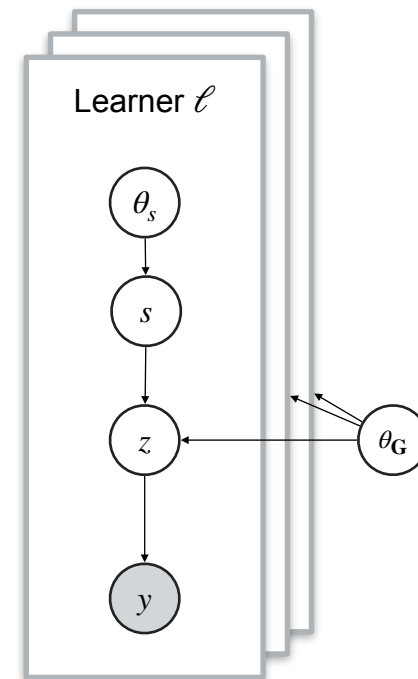
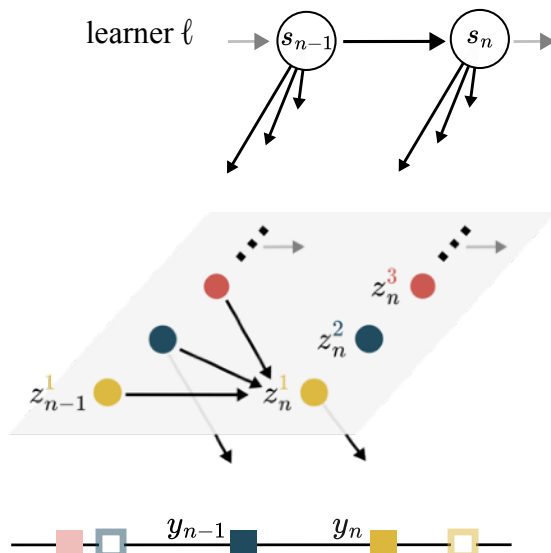
Latent  
cognitive traits  $s$

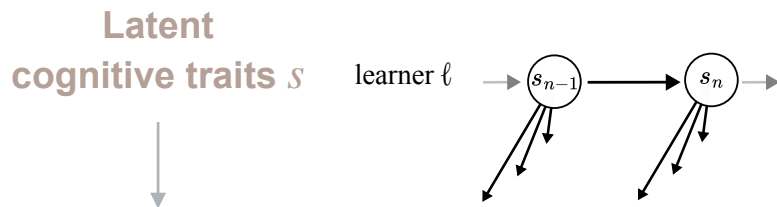


Latent  
knowledge states  $z$



Observed  
performance  $y$



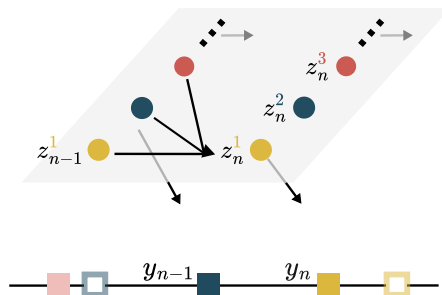


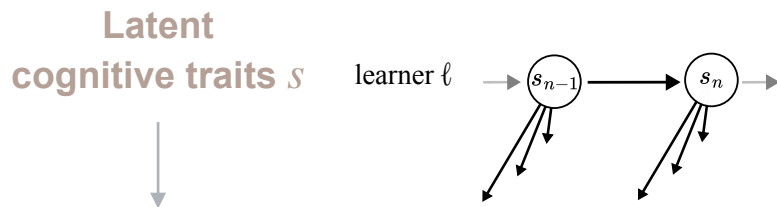
Per learner  $s_n^\ell := (\alpha_n^\ell, \mu_n^\ell, \gamma_n^\ell, \sigma_n^\ell)$  for personalization

- Memory: forgetting rate  $\alpha_n^\ell$ , long-term consolidation  $\mu_n^\ell$
- Structure: transfer ability  $\gamma_n^\ell$
- Noise: knowledge volatility  $\sigma_n^\ell$

Latent  
knowledge states  $z$

Observed  
performance  $y$



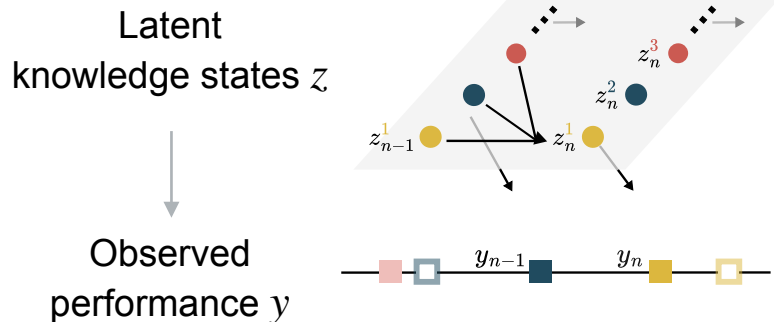


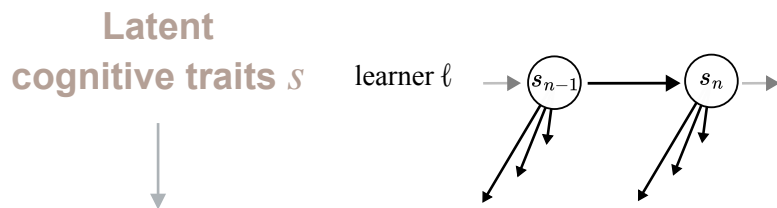
Per learner  $s_n^\ell := (\alpha_n^\ell, \mu_n^\ell, \gamma_n^\ell, \sigma_n^\ell)$  for personalization

- Memory: forgetting rate  $\alpha_n^\ell$ , long-term consolidation  $\mu_n^\ell$
- Structure: transfer ability  $\gamma_n^\ell$
- Noise: knowledge volatility  $\sigma_n^\ell$

Evolution:

- $s_n^\ell \sim p_{\theta_s}(s_n^\ell | s_{n-1}^\ell) := \mathcal{N}(s_n^\ell | Hs_{n-1}^\ell, R)$

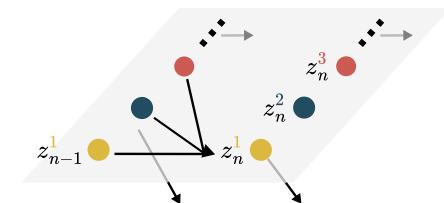




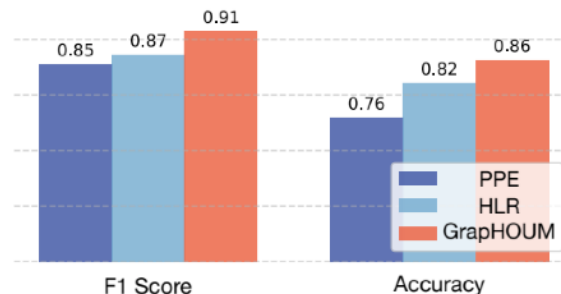
Per learner  $s_n^\ell := (\alpha_n^\ell, \mu_n^\ell, \gamma_n^\ell, \sigma_n^\ell)$  for personalization

- Memory: forgetting rate  $\alpha_n^\ell$ , long-term consolidation  $\mu_n^\ell$
- Structure: transfer ability  $\gamma_n^\ell$
- Noise: knowledge volatility  $\sigma_n^\ell$

Latent  
knowledge states  $z$



Observed  
performance  $y$



Zhou, H., Tejero-Cantero, A., & Wu, C. M (2023 CCN). The Dynamic and Structured Nature of Learning and Memory.

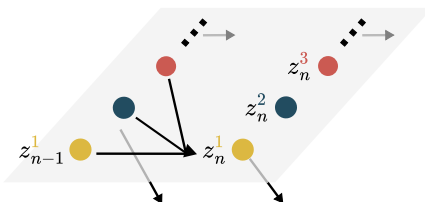
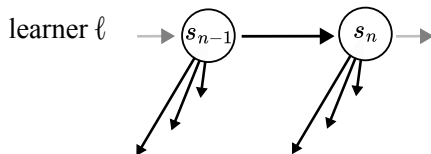
Latent  
cognitive traits  $s$



Latent  
knowledge states  $z$



Observed  
performance  $y$



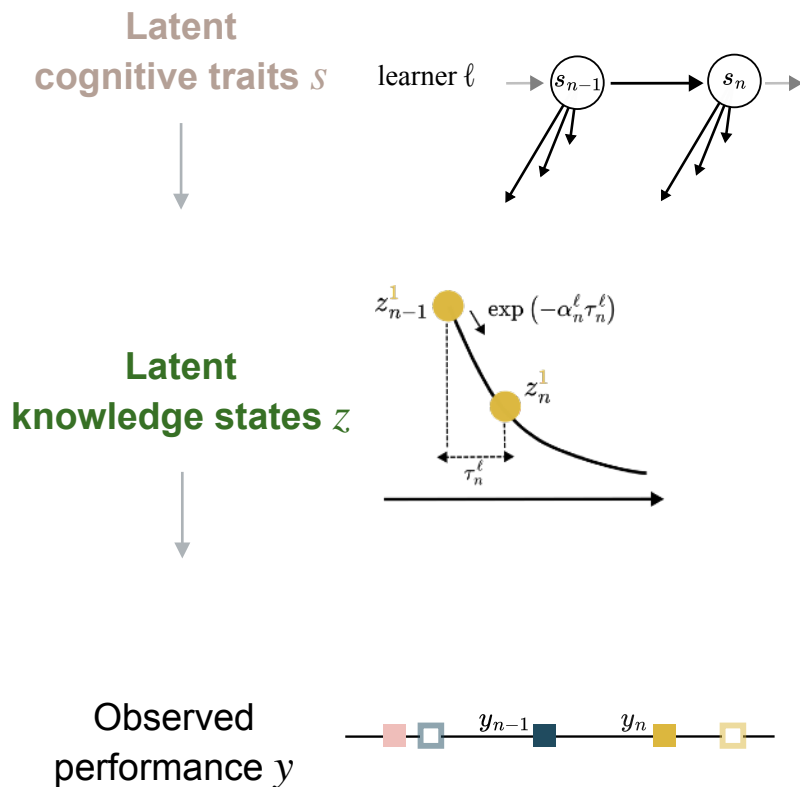
Per learner  $s_n^\ell := (\alpha_n^\ell, \mu_n^\ell, \gamma_n^\ell, \sigma_n^\ell)$  for personalization

- Memory: forgetting rate  $\alpha_n^\ell$ , long-term consolidation  $\mu_n^\ell$
- Structure: transfer ability  $\gamma_n^\ell$
- Noise: knowledge volatility  $\sigma_n^\ell$

Per learner & KC  $z_n^{\ell,k}$  for learning dynamics

Ornstein-Uhlenbeck process

$$dz^{\ell,k}/dt = \alpha^\ell (\mu^\ell - z^{\ell,k}) + \sigma^\ell \eta(t)$$



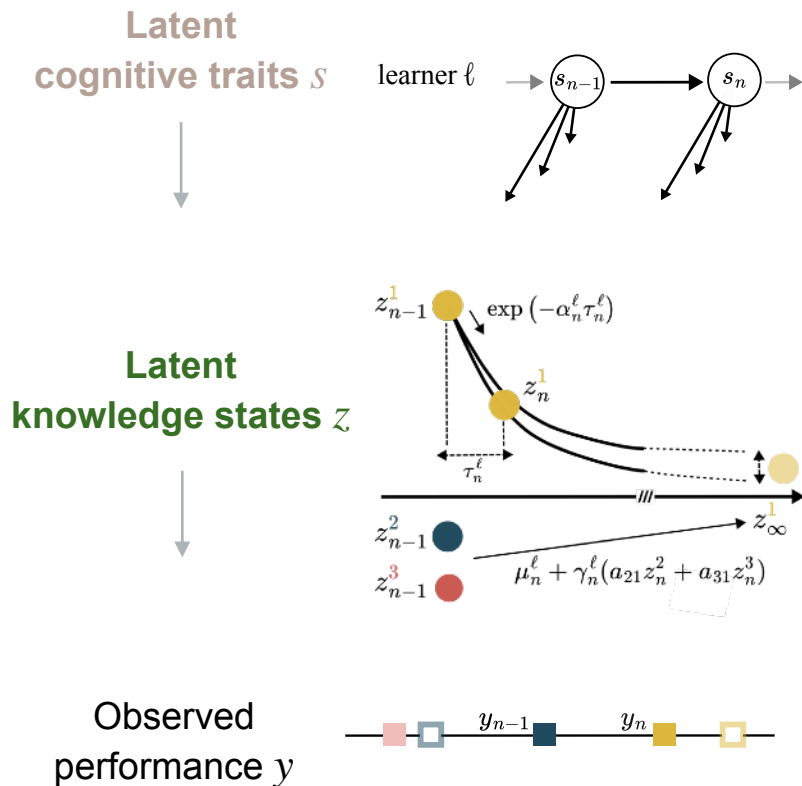
Per learner  $s_n^\ell := (\alpha_n^\ell, \mu_n^\ell, \gamma_n^\ell, \sigma_n^\ell)$  for personalization

- Memory: forgetting rate  $\alpha_n^\ell$ , long-term consolidation  $\mu_n^\ell$
- Structure: transfer ability  $\gamma_n^\ell$
- Noise: knowledge volatility  $\sigma_n^\ell$

Per learner & KC  $z_n^{\ell,k}$  for learning dynamics

$$dz^{\ell,k}/dt = \alpha^\ell (\mu^\ell - z^{\ell,k}) + \sigma^\ell \eta(t)$$

- Short-term: exponential decay
  - $z_n^{\ell,k} = z_{n-1}^{\ell,k} \exp(-\alpha_n^\ell \tau_n^\ell)$



Per learner  $s_n^\ell := (\alpha_n^\ell, \mu_n^\ell, \gamma_n^\ell, \sigma_n^\ell)$  for personalization

- Memory: forgetting rate  $\alpha_n^\ell$ , long-term consolidation  $\mu_n^\ell$
- Structure: transfer ability  $\gamma_n^\ell$
- Noise: knowledge volatility  $\sigma_n^\ell$

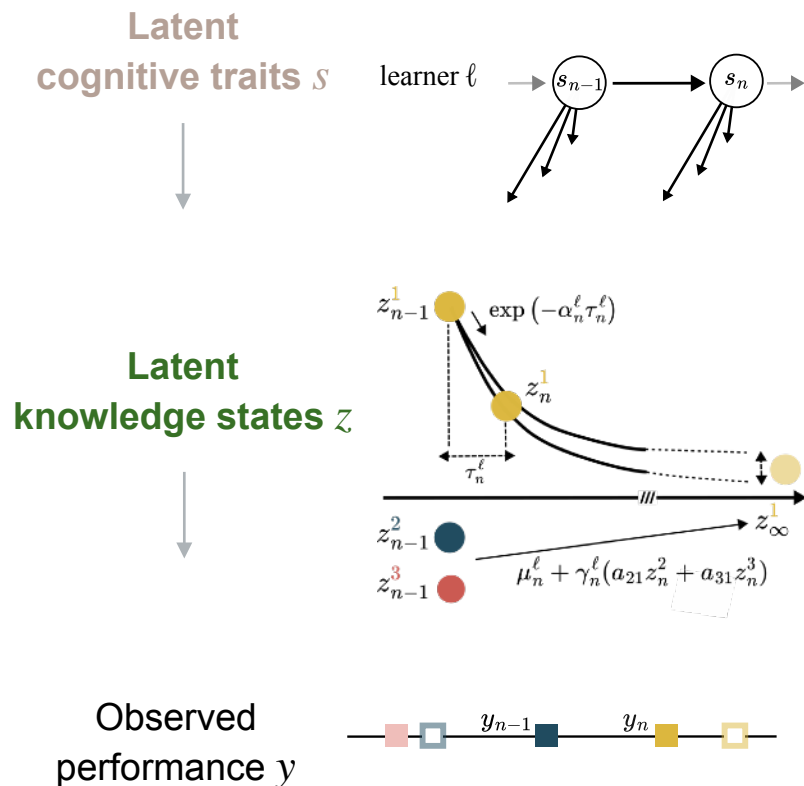
Per learner & KC  $z_n^{\ell,k}$  for learning dynamics

$$dz^{\ell,k}/dt = \alpha^\ell (\mu^\ell - z^{\ell,k}) + \sigma^\ell \eta(t)$$

- Short-term: exponential decay
  - $z_n^{\ell,k} = z_{n-1}^{\ell,k} \exp(-\alpha_n^\ell \tau_n^\ell)$
- Long-term: shifted by global structure

$$\tilde{\mu}_n^{\ell,k} := \mu_n^\ell + \gamma_n^\ell \sum_{i \neq k} a_{ik} z_n^{\ell,i}, \quad a_{ik} \in \theta_G$$





Per learner  $s_n^\ell := (\alpha_n^\ell, \mu_n^\ell, \gamma_n^\ell, \sigma_n^\ell)$  for personalization

- Memory: forgetting rate  $\alpha_n^\ell$ , long-term consolidation  $\mu_n^\ell$
- Structure: transfer ability  $\gamma_n^\ell$
- Noise: knowledge volatility  $\sigma_n^\ell$

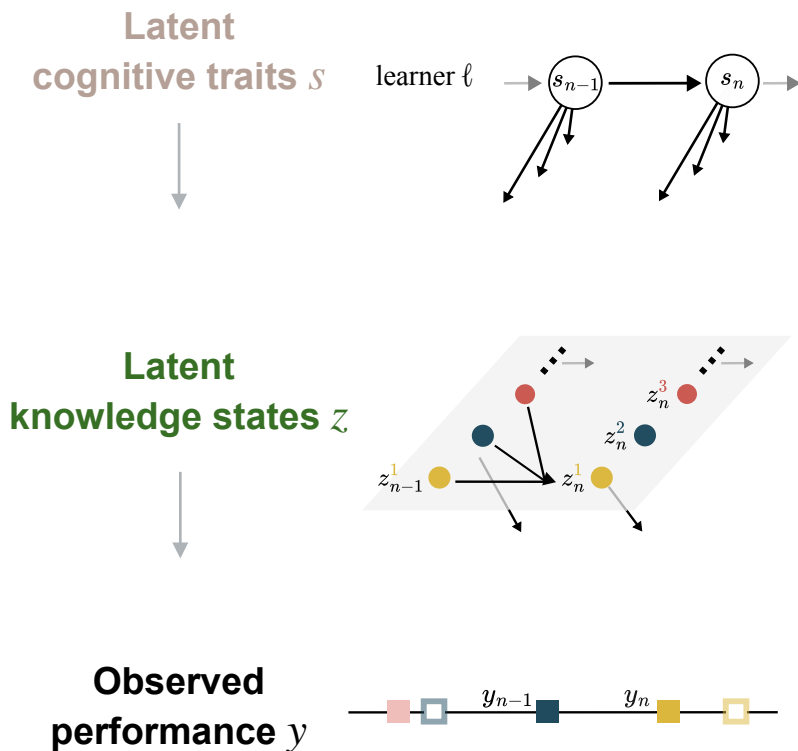
Per learner & KC  $z_n^{\ell,k}$  for learning dynamics

- Short-term: exponential decay
  - $dz^{\ell,k}/dt = \alpha^\ell (\mu^\ell - z^{\ell,k}) + \sigma^\ell \eta(t)$
- Long-term: shifted by global structure

$$\tilde{\mu}_n^{\ell,k} := \mu_n^\ell + \gamma_n^\ell \sum_{i \neq k} a_{ik} z_n^{\ell,i}, \quad a_{ik} \in \theta_G$$

Evolution:

$$z_n^{\ell,k} = \underbrace{z_{n-1}^{\ell,k} \exp(-\alpha_n^\ell \tau_n^\ell)}_{\text{short-term dynamics}} + \underbrace{\tilde{\mu}_n^{\ell,k} (1 - \exp(-\alpha_n^\ell \tau_n^\ell))}_{\text{long-term dynamics}}$$



**Per learner  $s_n^\ell := (\alpha_n^\ell, \mu_n^\ell, \gamma_n^\ell, \sigma_n^\ell)$  for personalization**

- Memory: forgetting rate  $\alpha_n^\ell$ , long-term consolidation  $\mu_n^\ell$
- Structure: transfer ability  $\gamma_n^\ell$
- Noise: knowledge volatility  $\sigma_n^\ell$

**Per learner & KC  $z_n^{\ell,k}$  for learning dynamics**

- Transient:  $dz^{\ell,k}/dt = \alpha^\ell (\mu^\ell - z^{\ell,k}) + \sigma^\ell \eta(t)$
- Long-term:  $\tilde{\mu}_n^{\ell,k^\dagger} := \mu_n^\ell + \gamma_n^\ell \sum_{i \neq k} a_{ik} z_n^{\ell,i}$

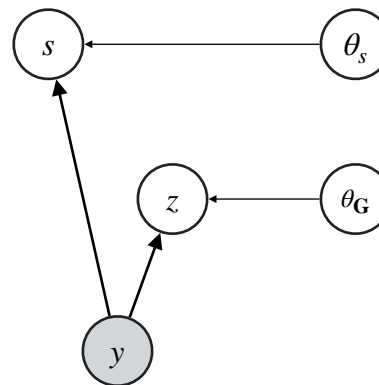
**Observation**

- Emission:  $\hat{y}_n^\ell \sim p(y_n^\ell | z_n^{\ell,k}) := \text{Bern}(\text{sigmoid}(z_n^{\ell,k}))$



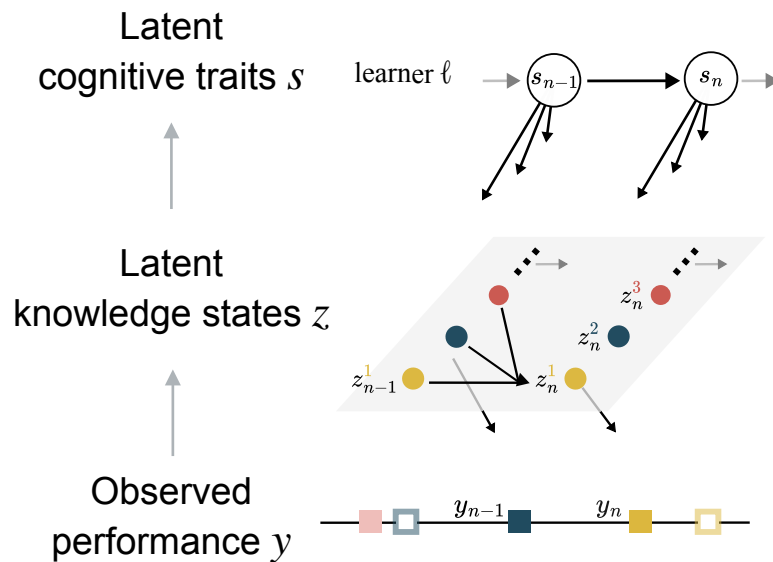
Exact inference over latent variables

- Full distribution over latents
- Point estimation over parameters



Exact inference over latent variables 🤖

$$p_{\theta}(s_{1:n}, \mathbf{z}_{1:n} \mid y_{1:n})$$
$$= \frac{p_{\theta}(s_{1:n}, \mathbf{z}_{1:n}, y_{1:n})}{\int_{s_{1:n}} \int_{\mathbf{z}_{1:n}} p_{\theta}(s_{1:n}, \mathbf{z}_{1:n}, y_{1:n})}$$

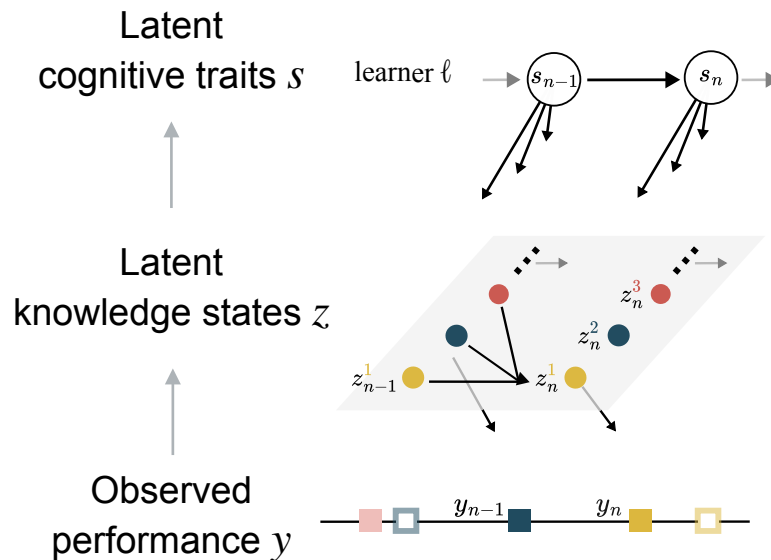


We need help from NN-based  
Approximate Bayesian Inference 🧐

$$q_{\phi}(z_{1:n}, s_{1:n} \mid y_{1:n}) = q_{\phi}(z_{1:n}) q_{\phi}(s_{1:n})$$

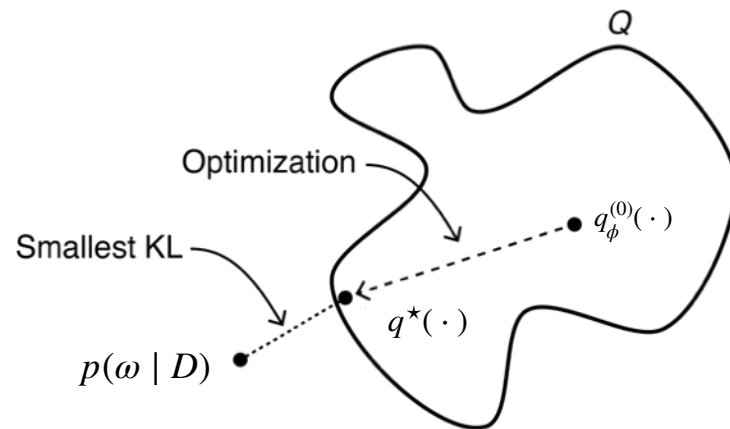


$$p_{\theta}(s_{1:n}, z_{1:n} \mid y_{1:n})$$



We need help from NN-based  
Approximate Bayesian Inference 🧐

$$q^*(\omega) \\ = \arg \min_{q(\cdot) \in \mathcal{Q}} \text{KL}[q_\phi(\omega) \| p_\theta(\omega | D)]$$



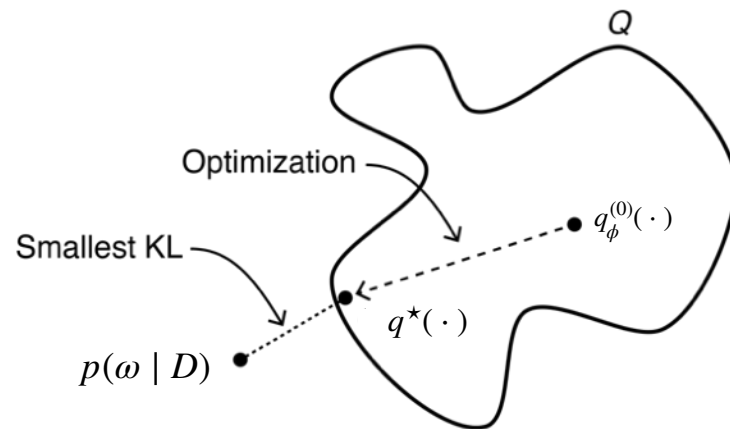
Latent states  $\omega := \{s_{1:n}, \mathbf{z}_{1:n}\}$

Observation  $D := \mathcal{H}_{1:n}$

Figure source: <https://gregorygundersen.com/blog/2021/04/16/variational-inference/>

We need help from NN-based  
Approximate Bayesian Inference 🧐

$$q^*(\omega) \\ = \arg \min_{q(\cdot) \in \mathcal{Q}} \text{KL}[q_\phi(\omega) \| p_\theta(\omega \mid D)]$$



$$q^*(\mathbf{z}_{1:n})q^*(s_{1:n}) \\ = \arg \min_{q(\cdot) \in \mathcal{Q}} \text{KL}[q_\phi(\mathbf{z}_{1:n})q_\phi(s_{1:n}) \| p_\theta(s_{1:n}, \mathbf{z}_{1:n} \mid y_{1:n})]$$

Latent states  $\omega := \{s_{1:n}, \mathbf{z}_{1:n}\}$

Observation  $D := \mathcal{H}_{1:n}$

Figure source: <https://gregorygundersen.com/blog/2021/04/16/variational-inference/>

- ELBO

$$\begin{aligned}\text{KL} \left[ q_{\phi}(\omega) \| p_{\theta}(\omega \mid D) \right] &= \int_{\omega} q_{\phi}(\omega) \log \frac{q_{\phi}(\omega)}{p_{\theta}(\omega \mid D)} = \mathbb{E}_q \left[ \log \frac{q_{\phi}(\omega)}{p_{\theta}(\omega \mid D)} \right] \\ &= \underbrace{-\mathbb{E}_q[\log p_{\theta}(D \mid \omega)] + \mathbb{E}_q[\log \frac{q_{\phi}(\omega)}{p_{\theta}(\omega)}]}_{-\text{ELBO}} + \log p(D)\end{aligned}$$

$$L_{\text{ELBO}}(\phi, \theta) = \mathbb{E}_{q_{\phi}(\omega)} [\log p_{\theta}(D \mid \omega)] - \text{KL} \left[ q_{\phi}(\omega) \| p_{\theta}(\omega) \right]$$

Latent states  $\omega := \{s_{1:n}, \mathbf{z}_{1:n}\}$

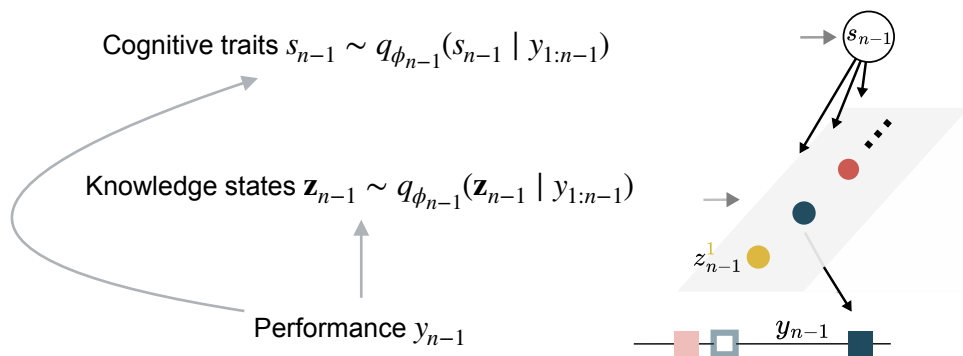
Observation  $D := \mathcal{H}_{1:n}$



- ELBO for fixed learning histories

$$\begin{aligned} L_{\text{ELBO}}(\phi, \theta) &= \mathbb{E}_{q_{\phi}(\omega)} \left[ \log p_{\theta}(D \mid \omega) \right] - \text{KL} \left[ q_{\phi}(\omega) \parallel p_{\theta}(\omega) \right] \\ &= \mathbb{E}_{q_{\phi}(z_{1:n}, s_{1:n})} \left[ \log p_{\theta}(y_{1:n} \mid z_{1:n}, s_{1:n}) - \log(q_{\phi}(z_{1:n}, s_{1:n}) - p_{\theta}(z_{1:n}, s_{1:n} \mid z_0, s_0)) \right] \end{aligned}$$

- ELBO for online data

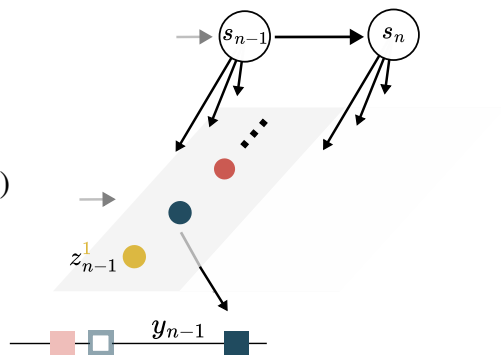


- ELBO for online data

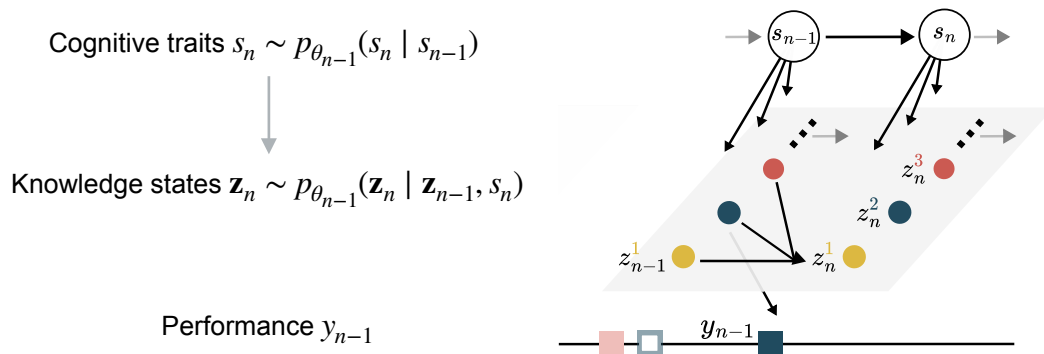
Cognitive traits  $s_{n-1} \sim q_{\phi_{n-1}}(s_{n-1} \mid y_{1:n-1})$

Knowledge states  $\mathbf{z}_{n-1} \sim q_{\phi_{n-1}}(\mathbf{z}_{n-1} \mid y_{1:n-1})$

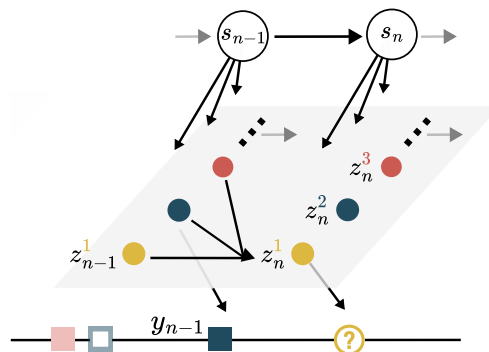
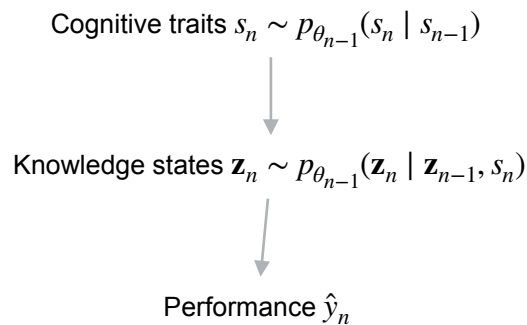
Performance  $y_{n-1}$



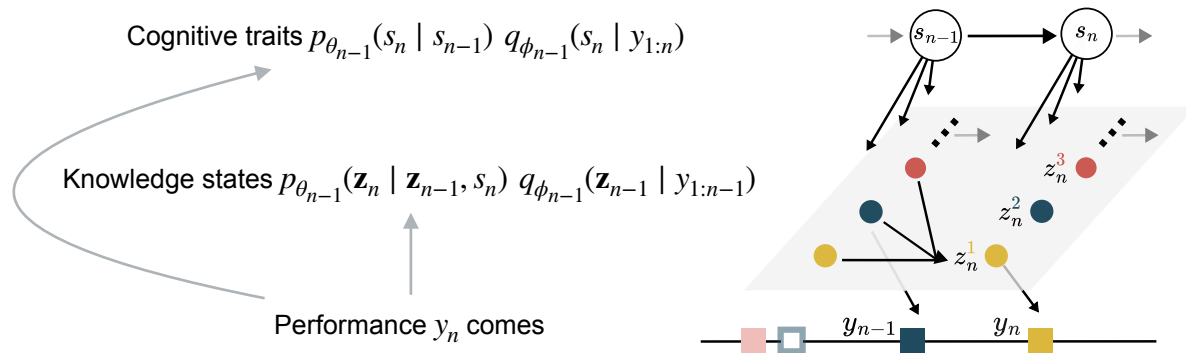
- ELBO for online data



- ELBO for online data



- ELBO for online data



- ELBO for online data

$$L_{\text{ELBO}}(\phi, \theta) = \mathbb{E}_{q_{\phi}(\omega_{1:n})} \left[ \log p_{\theta}(D_{1:n} \mid \omega_{1:n}) \right] - \text{KL} \left[ q_{\phi}(\omega_{1:n}) \parallel p_{\theta}(\omega_{1:n}) \right]$$

$$L_{\text{ELBO}}(\phi_n, \theta_n) = \mathbb{E}_{q_{\phi_n}(\omega_n)} \left[ \log p_{\theta_n}(D_n \mid \omega_n) \right] - \text{KL} \left[ q_{\phi_n}(\omega_n) \parallel q_{\phi_{n-1}}(\omega_{n-1}) \right]$$



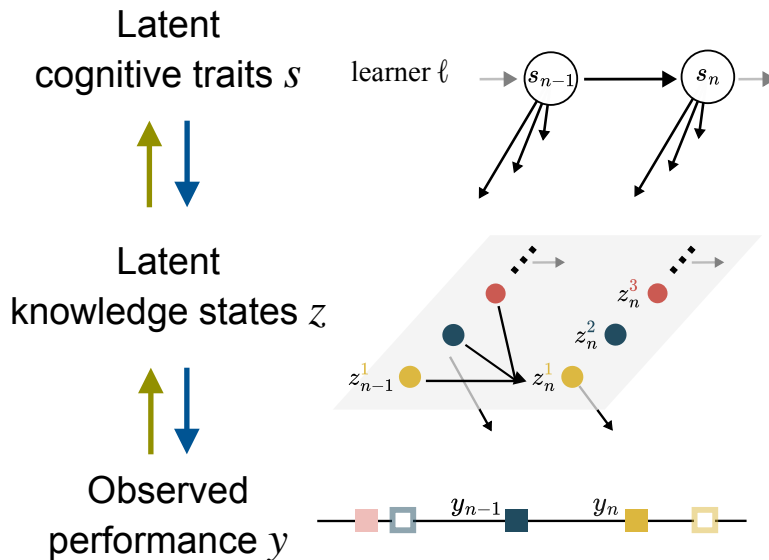
## Latent inference:

Online variational inference

$$L_{\text{ELBO}}(\phi, \theta)$$

$$= \mathbb{E}_{q_{\phi}(\omega)} [\log p_{\theta}(D \mid \omega)]$$

$$- \text{KL} [q_{\phi}(\omega) \parallel p_{\theta}(\omega)]$$



Per learner  $s_n^{\ell} := (\alpha_n^{\ell}, \mu_n^{\ell}, \gamma_n^{\ell}, \sigma_n^{\ell})$

- forgetting rate  $\alpha_n^{\ell}$ , long-term consolidation  $\mu_n^{\ell}$
- transfer ability  $\gamma_n^{\ell}$
- knowledge volatility  $\sigma_n^{\ell}$

Per learner & KC  $z_n^{\ell, k}$

- $dz^{\ell, k}/dt = \alpha^{\ell} (\mu^{\ell} - z^{\ell, k}) + \sigma^{\ell} \eta(t)$
- $\tilde{\mu}_n^{\ell, k \dagger} := \mu_n^{\ell} + \gamma_n^{\ell} \sum_{i \neq k} a_{ik} z_n^{\ell, i}$

## Observation

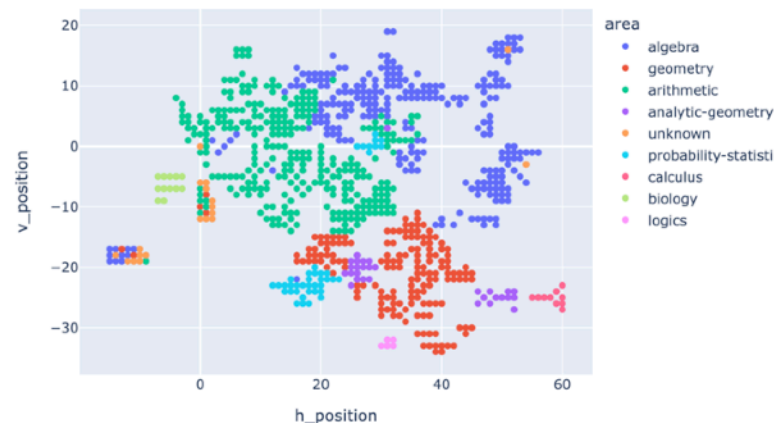
- Emission:  
 $\hat{y}_n^{\ell} \sim p(y_n^{\ell} \mid z_n^{\ell, k}) := \text{Bern}(\text{sigmoid}(z_n^{\ell, k}))$





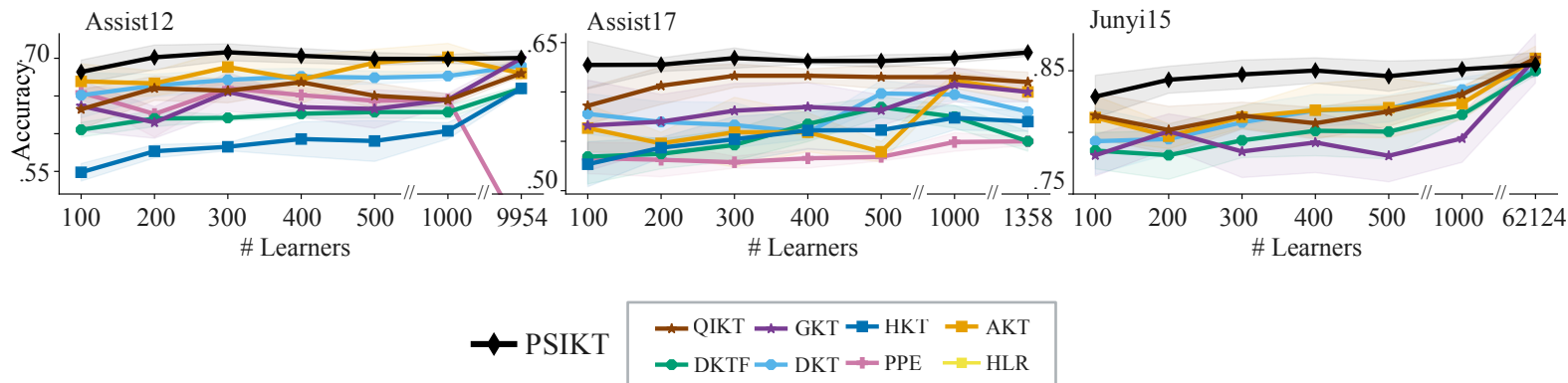
- Assistment 2012
- Assistment 2017
- Junyi 2015
  - Pre-college mathematics study

Exercises distribution on area in knowledge map

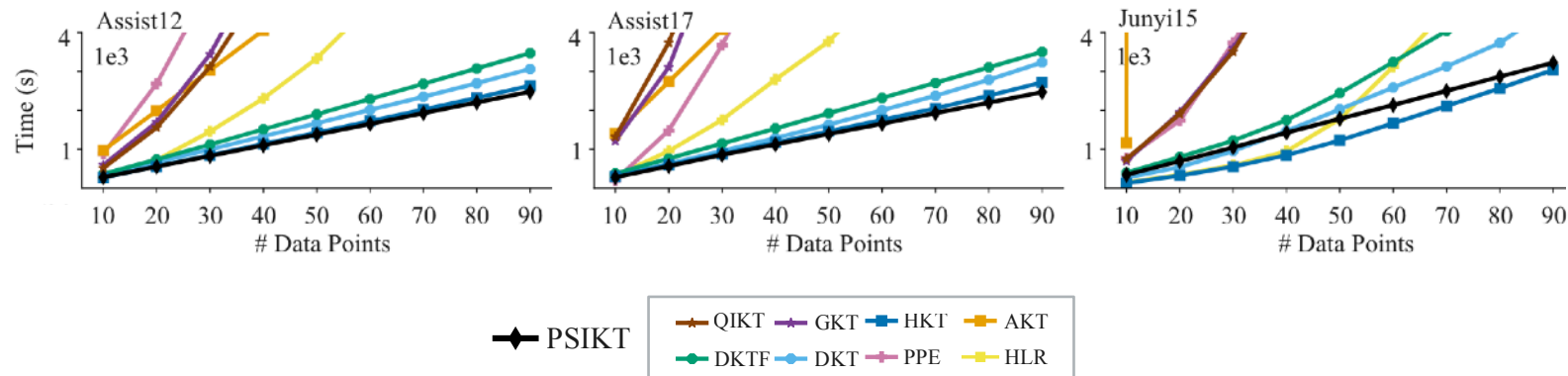


	user_id	exercise	problem_type	problem_number	topic_mode	suggested	review_mode	time_done	time_taken	time_taken_attempts	correct
0	12884	time_terminology	analog_word	1	False	False	False	1420714810324490	4	3&1	False
1	239464	multiplication_1	0	6	False	False	False	1403098400836660	2	2	True
2	147359	adding_decimals_0.5	0	6	False	False	False	1418890695540340	16	16	True
3	158155	multiplication_1	0	3	False	False	False	1400469444264040	2	2	True
4	147151	subtraction_2	subtraction-2	10	True	True	False	1382650905730160	4	4	True

- Within-learner 10-step prediction performance as a function of cohort sizes

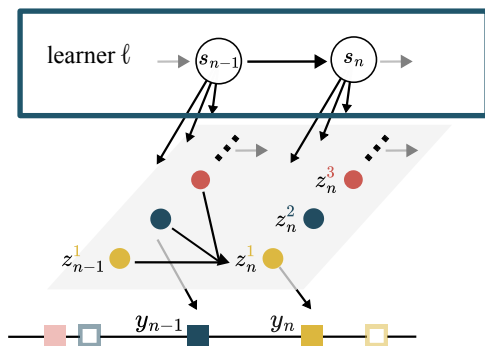


- Cumulative training time of continual learning





- **Cognitive traits:** forgetting rate, consolidation memory, transfer ability, volatility

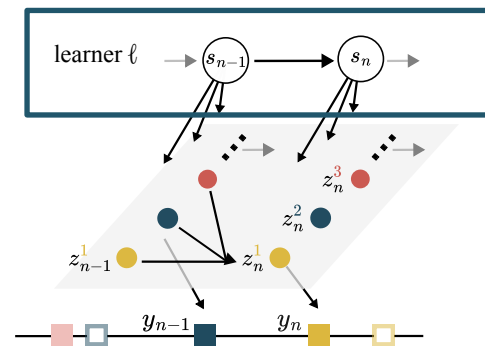




- **Cognitive traits:** forgetting rate, consolidation memory, transfer ability, volatility

## Representation capacity

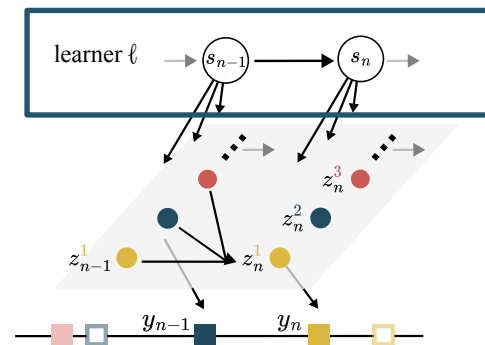
- *Specific* to each learner
- *Consistent* across data splits
- *Disentangled* across dimensions



- **Cognitive traits:** forgetting rate, consolidation memory, transfer ability, volatility

## Representation capacity

- *Specific* to each learner:  $MI(s; \ell) = H(s) - H(s \mid \ell)$
- *Consistent* across data splits
- *Disentangled* across dimensions



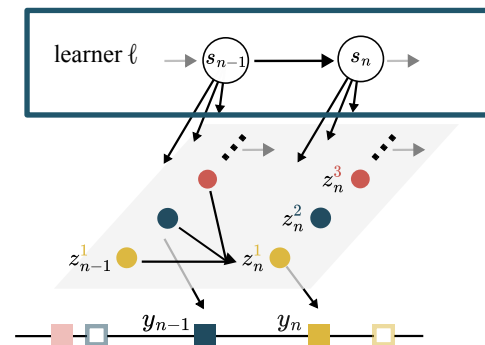
- **Cognitive traits:** forgetting rate, consolidation memory, transfer ability, volatility

## Representation capacity

- *Specific* to each learner
- *Consistent* across data splits:

$$\mathbb{E}_{\ell_{\text{sub}}} \text{MI}(s^\ell; \ell_{\text{sub}}) := \mathbb{E}_{\ell_{\text{sub}}} \left[ H(s \mid \ell) - H(s \mid \ell_{\text{sub}}) \right]$$

- *Disentangled* across dimensions

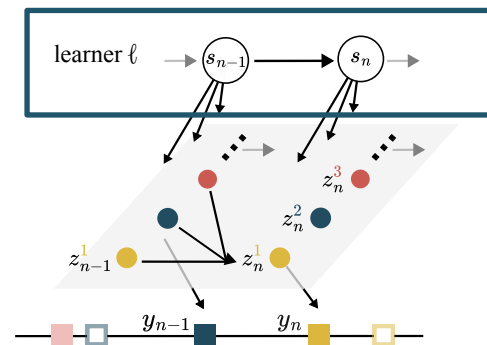


- **Cognitive traits:** forgetting rate, consolidation memory, transfer ability, volatility

## Representation capacity

- *Specific*
- *Consistent* across data splits
- *Disentangled* across dimensions:

$$\text{KL}(s||\ell) := H(s)_{\text{full}} - H(s | \ell)_{\text{diag}}$$





- **Cognitive traits:** forgetting rate, consolidation memory, transfer ability, volatility

## Representation capacity

- *Specific* to each learner
- *Consistent* across data splits
- *Disentangled* across dimensions

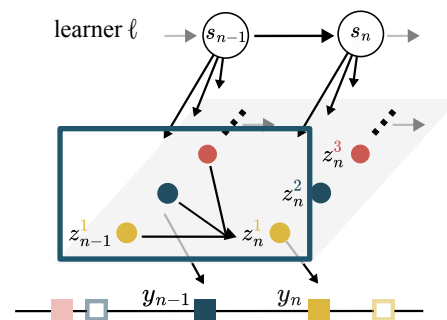
Metric	Dataset	Baseline	PSI-KT
Specificity $MI(s; \ell) \uparrow$	Assist12	<b>8.8</b>	<u>8.4</u>
	Assist17	<b>10.1</b>	<u>10.0</u>
	Junyi15	<u>13.5</u>	<b>14.4</b>
Consistency <sup>-1</sup> $\mathbb{E}_{\ell_{\text{sub}}} MI(s^{\ell}; \ell_{\text{sub}}) \downarrow$	Assist12	<u>12.2</u>	<b>7.4</b>
	Assist17	<b>6.4</b>	<b>6.4</b>
	Junyi15	<u>7.7</u>	<b>5.0</b>
Disentanglement $D_{KL}(s    \ell) \uparrow$	Assist12	<u>2.3</u>	<b>7.4</b>
	Assist17	<u>0.6</u>	<b>8.4</b>
	Junyi15	<u>5.0</u>	<b>11.5</b>

Bold indicates the better model.  
PSI-KT vs. the best baseline model.

- **Knowledge structure: prerequisites**

**Structure correctness:**

- Human-annotated ground-truth
- Learners' progress

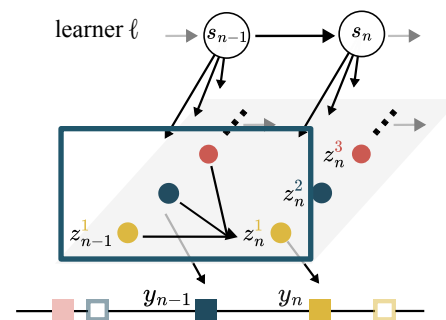


- **Knowledge structure: prerequisites**

## Structure correctness:

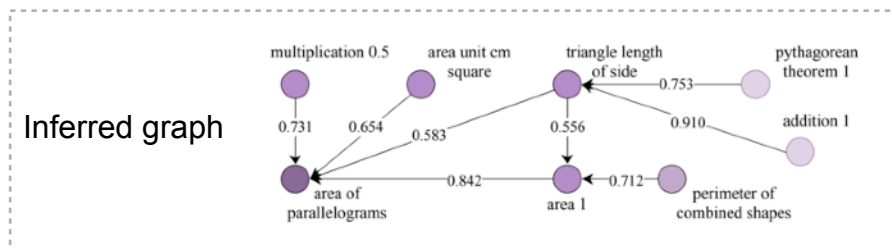
- Human-annotated ground-truth
- Learners' progress

Metric	MRR $\uparrow$	JS expert $\uparrow$	JS crowd $\uparrow$	nLL $\downarrow$
Dataset	Junyi15			
Best Baseline	.0082	.0015	.0047	<b>3.03</b>
PSI-KT	<b>.0086</b>	<b>.0019</b>	<b>.0095</b>	4.11

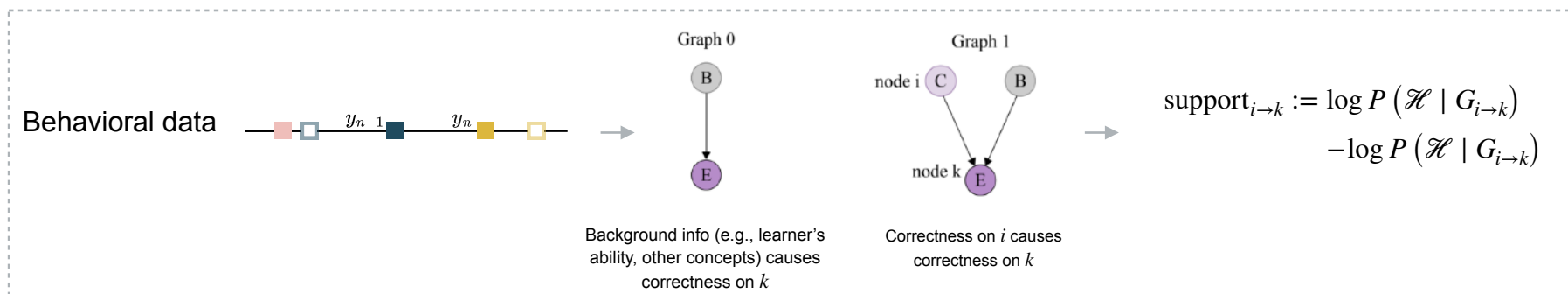


- **Knowledge structure:** prerequisites

## Structure correctness: causally support in human learning

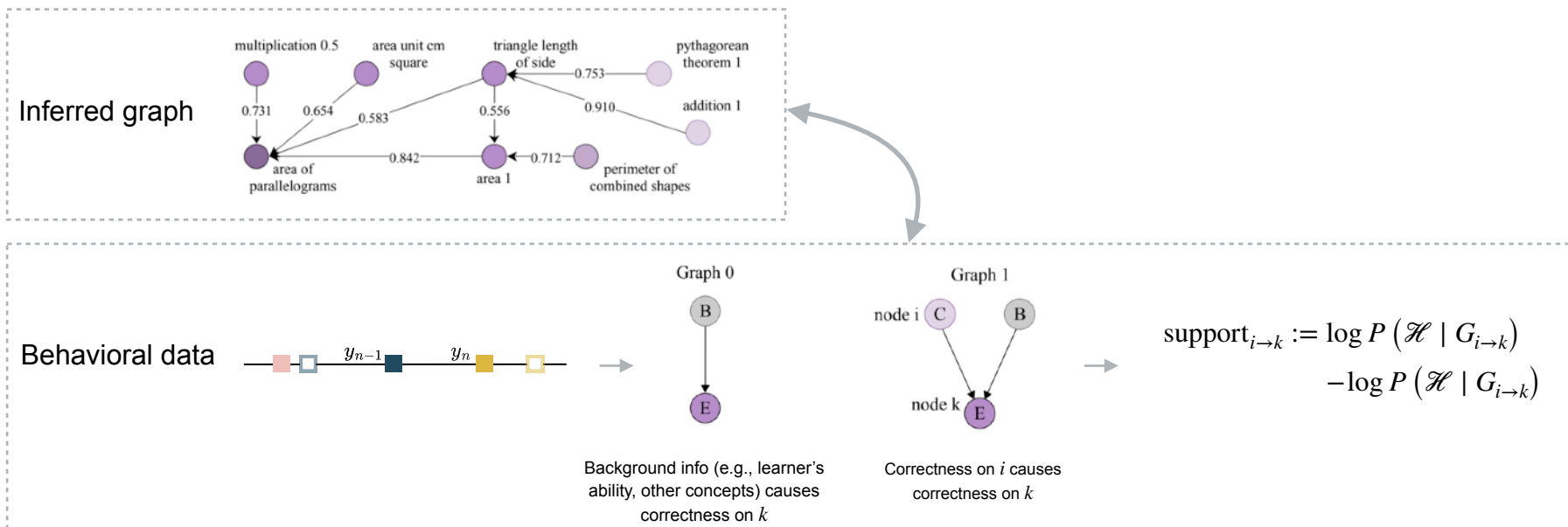


Griffiths, T. L., & Tenenbaum, J. B. (2009). Theory-based causal induction. *Psychological review*, 116(4), 661.



- **Knowledge structure: prerequisites**

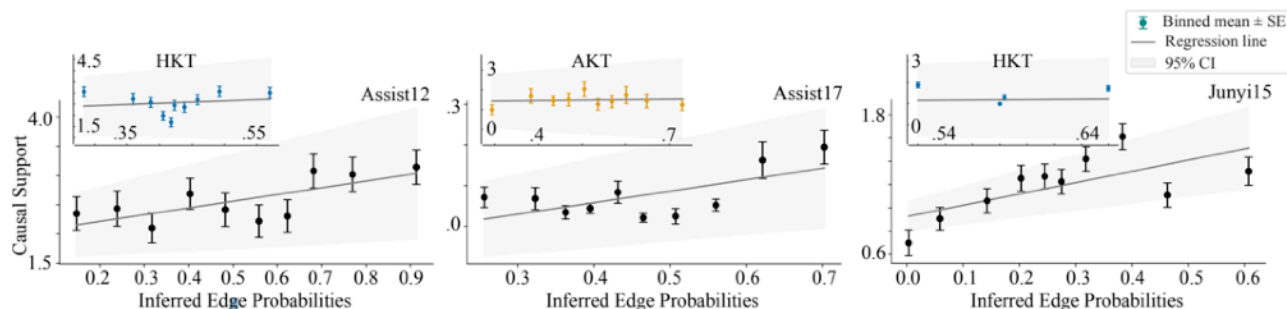
## Structure correctness: causally support in human learning



- **Knowledge structure:** prerequisites

**Structure correctness:** causally support in human learning

$$\text{support}_{i \rightarrow k} := \log P(\mathcal{H} \mid G_{i \rightarrow k}) - \log P(\mathcal{H} \mid G_{i \rightarrow k})$$

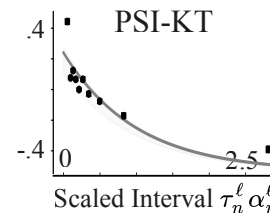
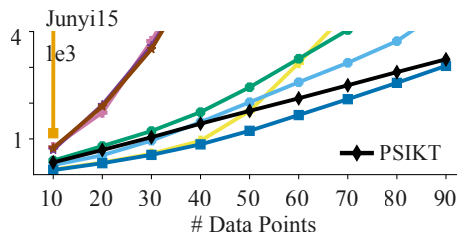
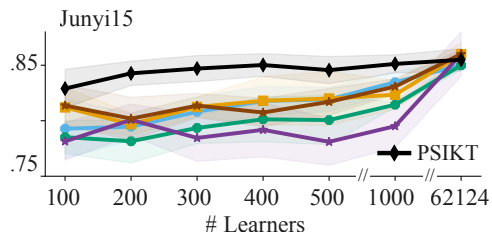


## PSI-KT

Predictive,

Scalable,

Interpretable



Knowledge tracing for future intelligent tutoring systems



Scan Me

# Acknowledgement



Robert Bamler



Charley Wu



Álvaro Tejero-Cantero

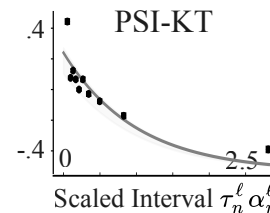
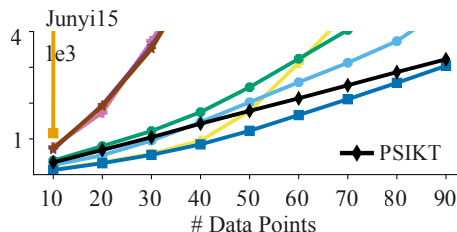
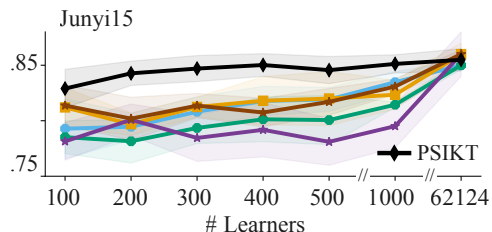


## PSI-KT

Predictive,

Scalable,

Interpretable



Knowledge tracing for future intelligent tutoring systems



Scan Me