

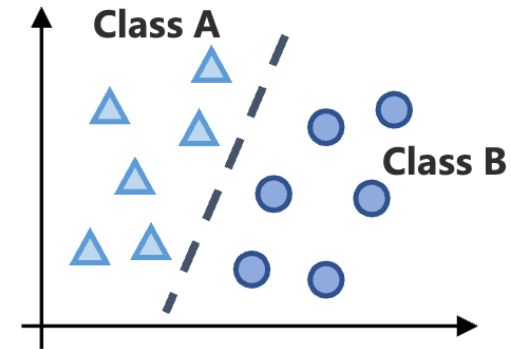
# TabICL: A Tabular Foundation Model for In-Context Learning on Large Data

Jingang QU, David Holzmüller, Gaël Varoquaux, Marine Le Morvan

ICML 2025

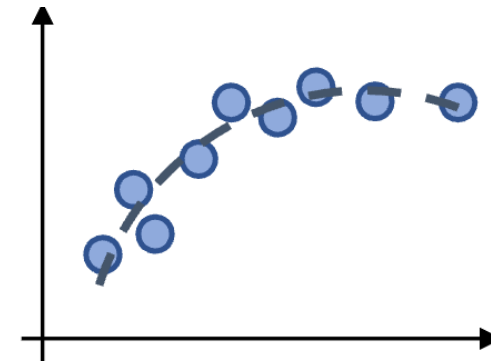
Currency	Amount	Card type	Age	<i>Fraud</i>
USD	3497.74	Debit	55	Yes
EUR	1121.53	Prepaid	45	No
CNY	2867.57	Credit	31	No
USD	4100.37	Debit	59	?

Classification



- Heterogenous features (numerical, categorical, ...)
- Missing values, uninformative features, and outliers
- Lack of spatial or sequential relationships

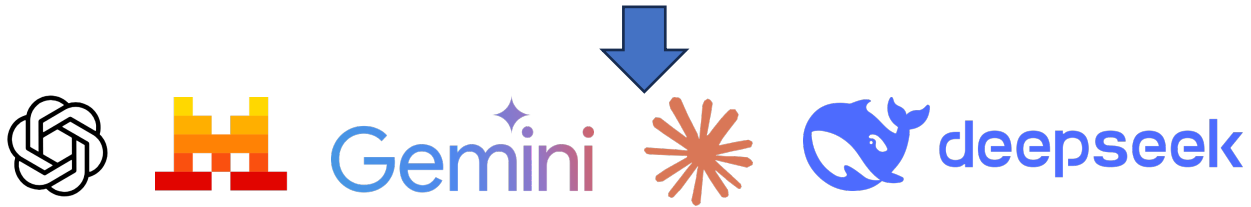
Regression




In-context examples

$7 + 2 = 9$	gaot => goat
$1 + 0 = 1$	sanke => snake
$3 + 4 = 7$	brid => bird
$5 + 9 = 14$	fsih => fish
$8 + 4 = 12$	dcuk => duck

Query  $9 + 8 = ?$   $cmihp => ?$



 The pattern here involves rearranging the jumbled letters to form the correct name of an animal. So,  $cmihp \Rightarrow chimp$  🐒

$\mathcal{D}_{\text{train}}$

$X \in \mathbb{R}^{n \times m}$        $y \in \mathbb{R}^{n \times 1}$

Columns (Features)

$x_{11}$	$x_{12}$	$\dots$	$x_{1m}$
$x_{21}$	$x_{22}$	$\dots$	$x_{2m}$
$\dots$	$\dots$	$\dots$	$\dots$
$x_{n1}$	$x_{n2}$	$\dots$	$x_{nm}$

Rows (Samples)

Target

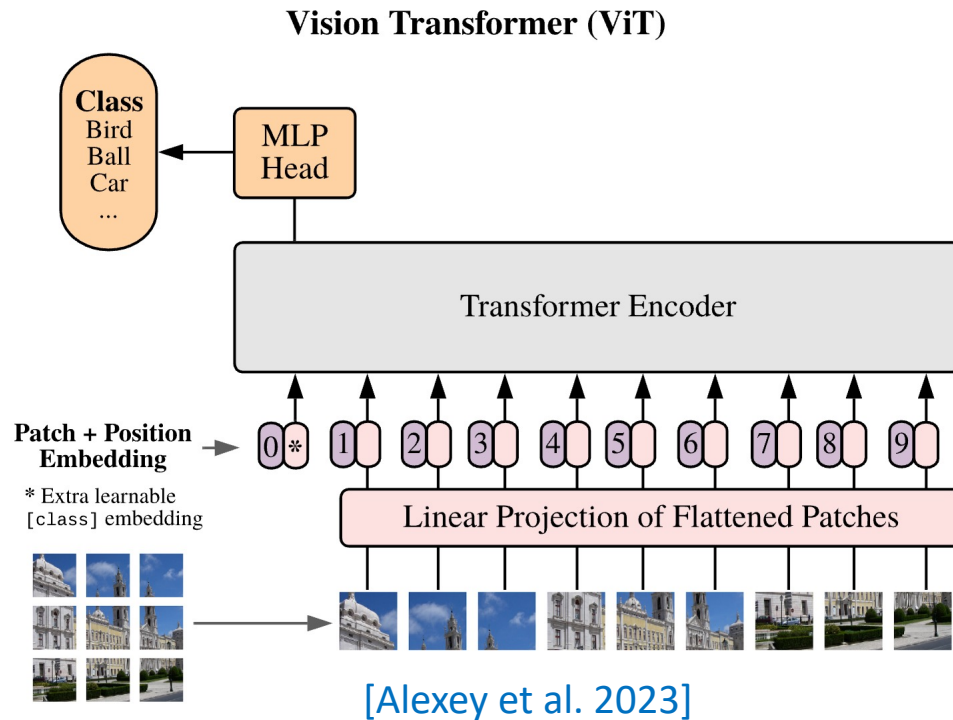
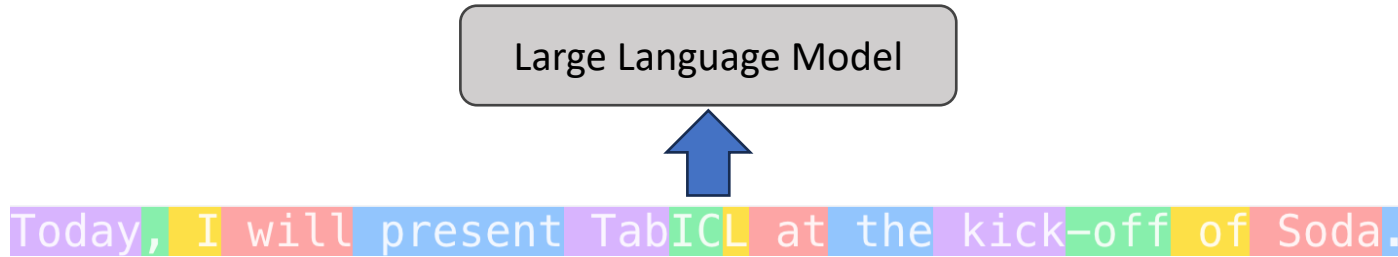
$y_1$
$y_2$
$\dots$
$y_n$

=>

Test  $x_1^* \ x_2^* \ \dots \ x_m^* \Rightarrow ?$

Transformer

$p(y^* | x^*, \mathcal{D}_{\text{train}}; \theta)$



$$X \in \mathbb{R}^{n \times m}$$

$$y \in \mathbb{R}^{n \times 1}$$

**Columns (Features)**

**Rows (Samples)**

$x_{11}$	$x_{12}$	$\dots$	$x_{1m}$
$x_{21}$	$x_{22}$	$\dots$	$x_{2m}$
$\dots$	$\dots$	$\dots$	$\dots$
$x_{n1}$	$x_{n2}$	$\dots$	$x_{nm}$

**Target**

$y_1$
$y_2$
$\dots$
$y_n$

$$\{(x_1, y_1), \dots, (x_n, y_n)\}, x^*$$



[Hollmann et al. 2023]

$$X \in \mathbb{R}^{n \times m}$$

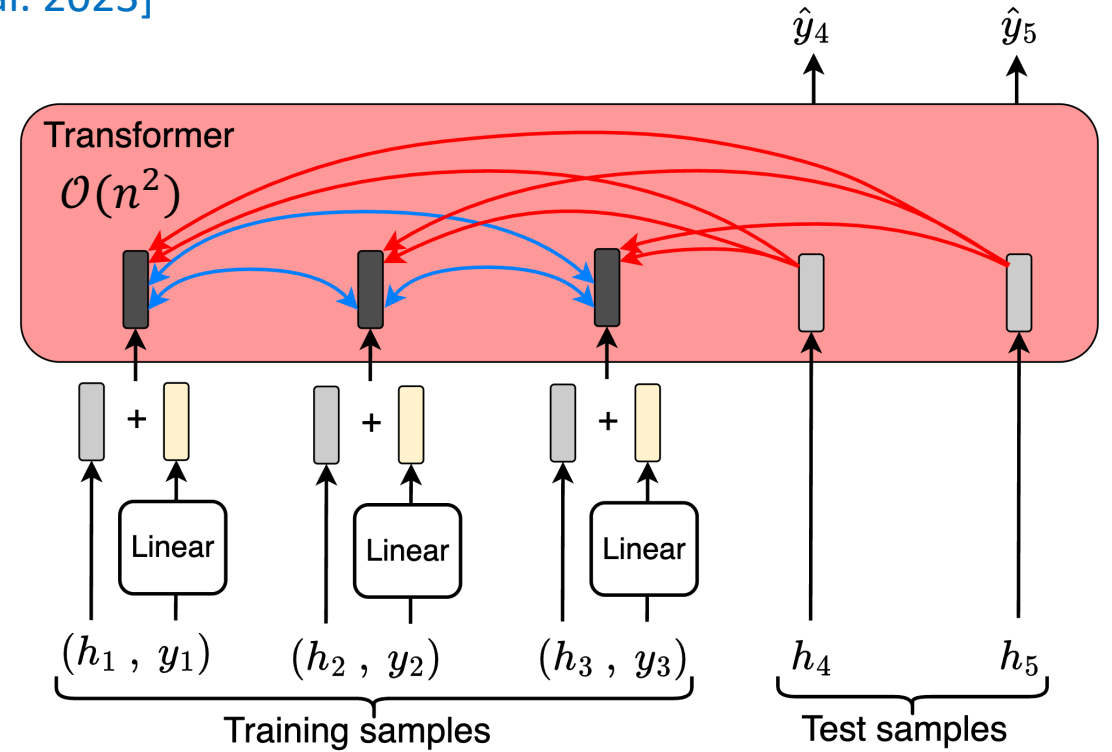
Columns (Features)

Rows (Samples)

$x_{11}$	$x_{12}$	$\dots$	$x_{1m}$
$x_{21}$	$x_{22}$	$\dots$	$x_{2m}$
$\dots$	$\dots$	$\dots$	$\dots$
$x_{n1}$	$x_{n2}$	$\dots$	$x_{nm}$

- Each row is a token
- $h_i = \text{Linear}(\text{row}_i)$
- Attention between  $h_i$  with  $\mathcal{O}(n^2)$

Complexity  $\Rightarrow \mathcal{O}(n^2)$



	Pretraining			Inference		
	# of samples	# of features	# of classes	# of samples	# of features	# of classes
TabPFN	1,024	100	10	3,000	100	10

[Hollmann et al. 2025]

$x_1$	$x_2$	$y$
1.2	6.1	3.0
8.9	9.1	3.1
1.0	2.9	6.7
33.3	2.2	?

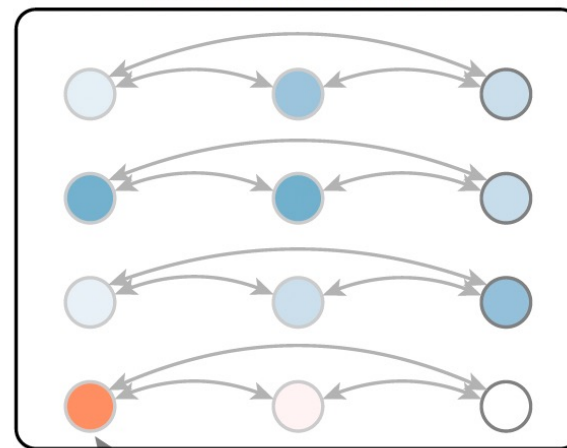
We predict this entry

Each cell is a token



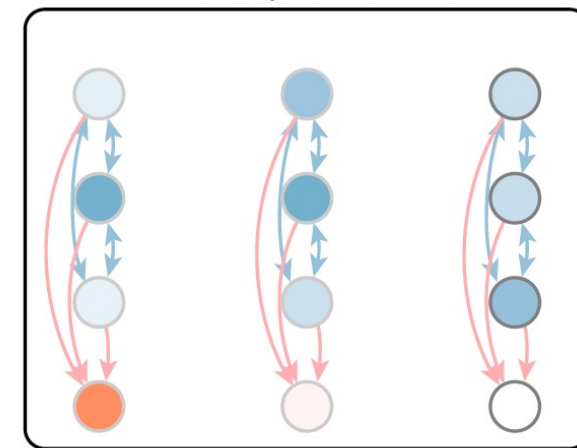
$$\mathcal{O}(n \cdot m^2)$$

1D feature attention



$$\mathcal{O}(m \cdot n^2)$$

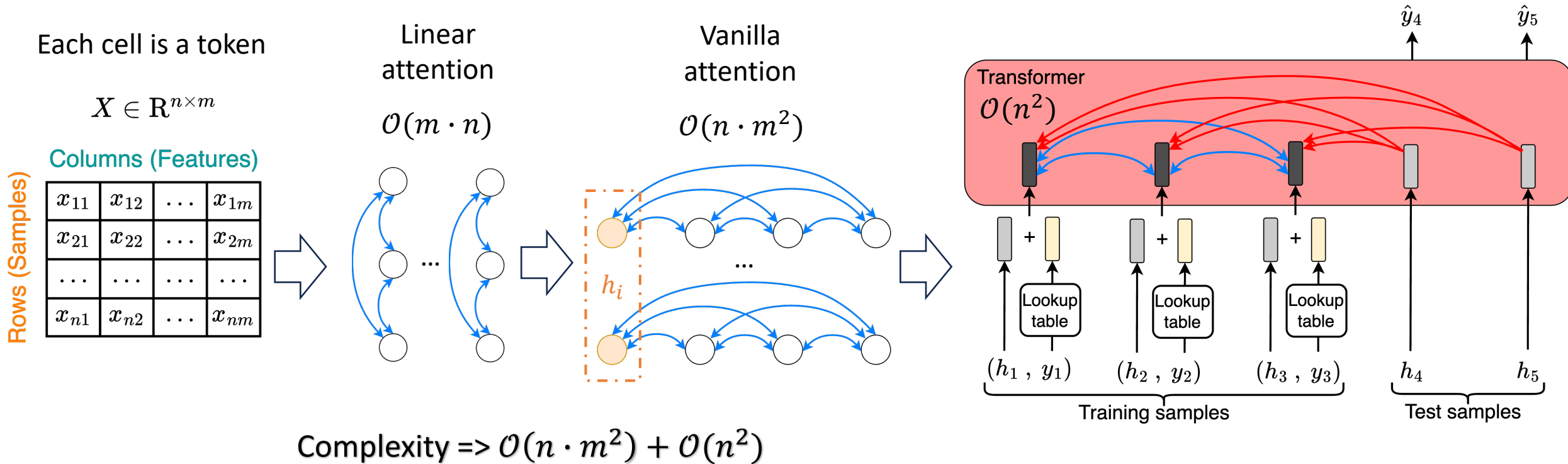
1D sample attention



Each node represents one entry in the table

$$\text{Complexity} \Rightarrow \mathcal{O}(n \cdot m^2) + \mathcal{O}(m \cdot n^2)$$

	<i>Pretraining</i>			<i>Inference</i>		
	# of samples	# of features	# of classes	# of samples	# of features	# of classes
TabPFN	1,024	100	10	3,000	100	10
TabPFNv2	2,048	160	10	10,000	500	10

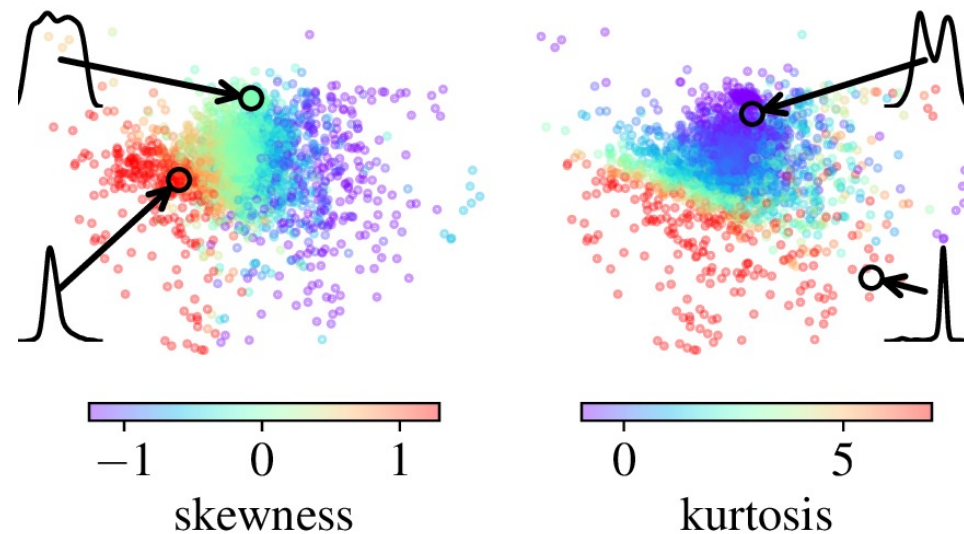
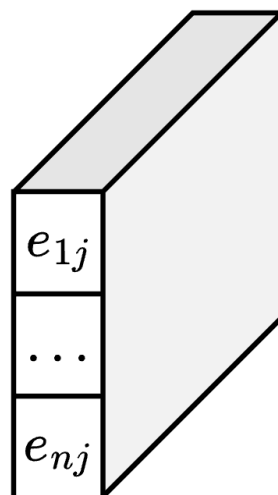
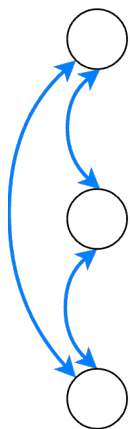
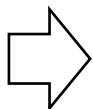
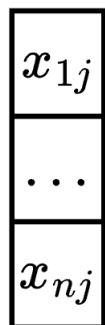


	Pretraining			Inference		
	# of samples	# of features	# of classes	# of samples	# of features	# of classes
TabPFN	1,024	100	10	3,000	100	10
TabPFNV2	2,048	160	10	10,000	500	10
TabICL	60,000	100	10	100,000	500	Any

Set Transformer  
(Linear attention)

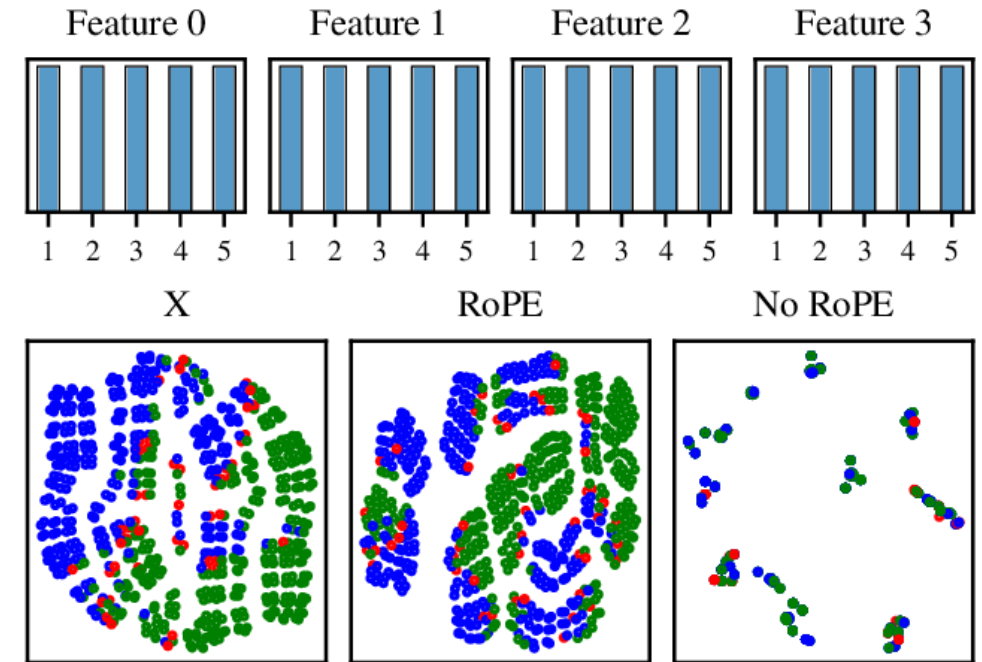
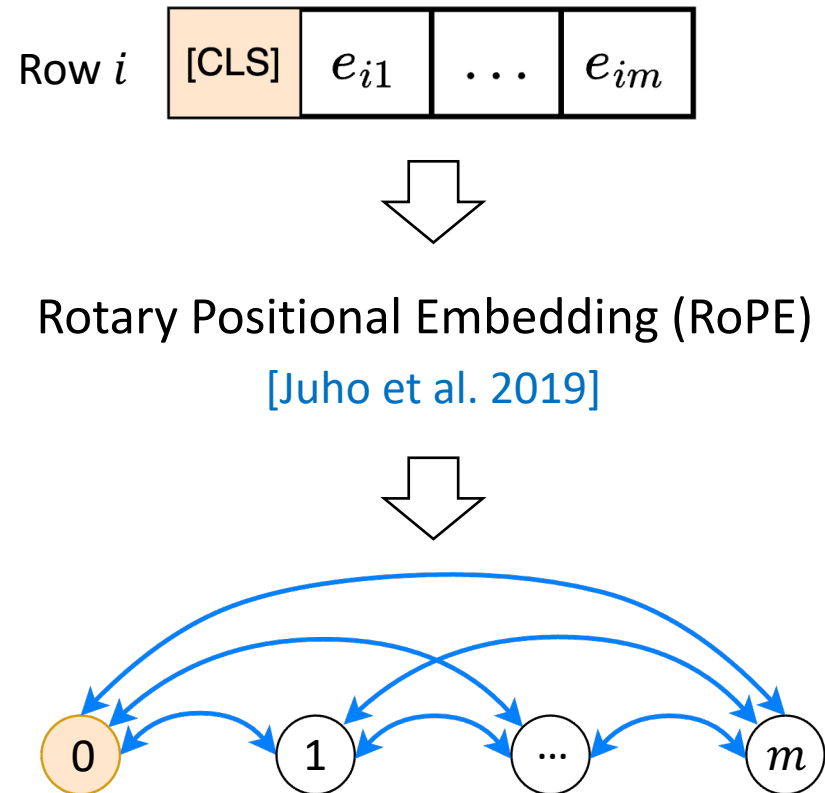
[Juho et al. 2019]

Column  $j$

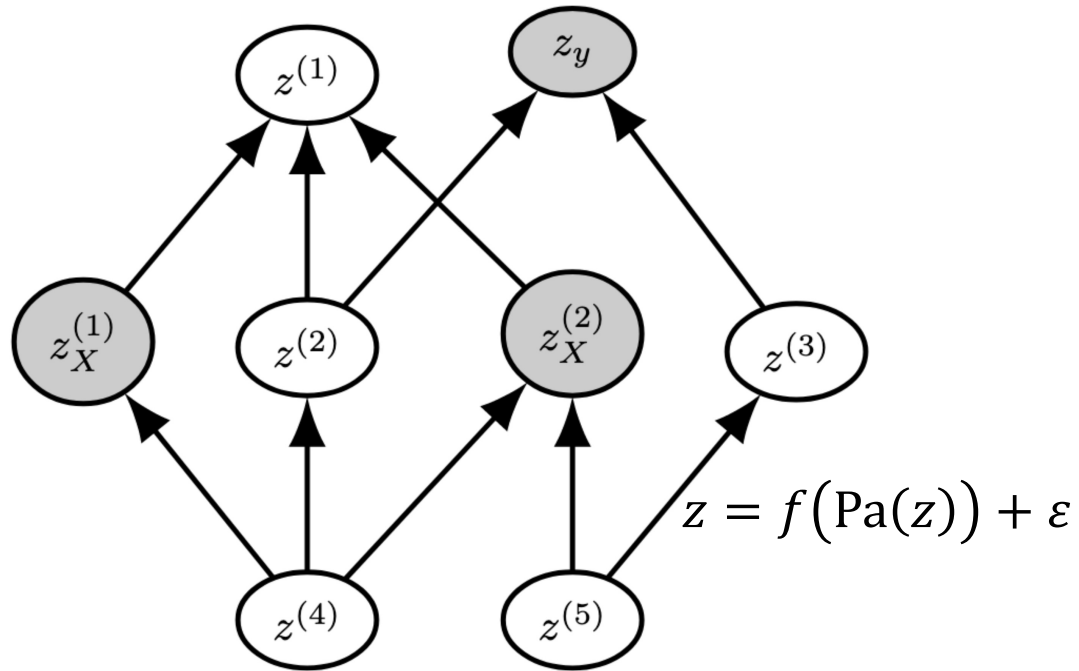


Learned embeddings encode statistical properties of feature distributions.





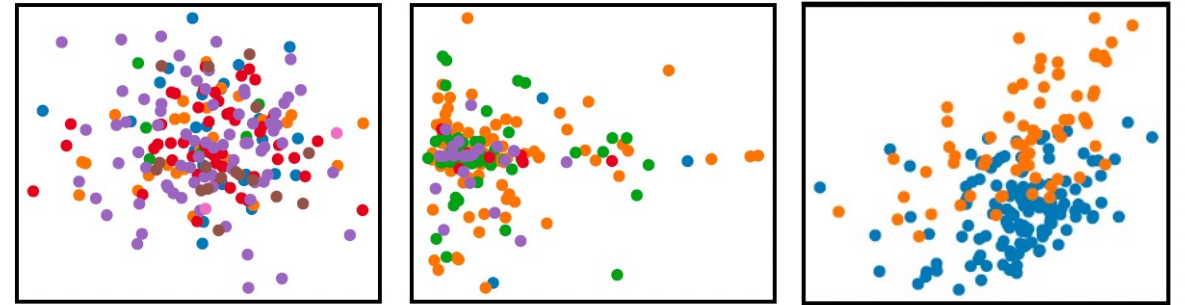
RoPE alleviates the representation collapse.



Structural causal models

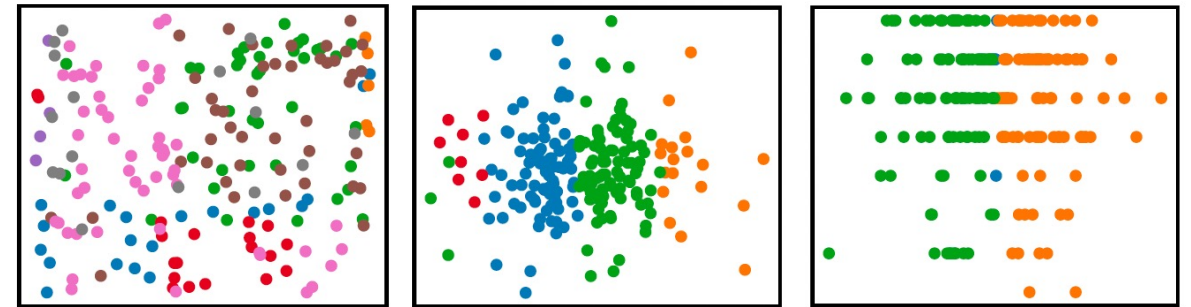
[Hollmann et al. 2022]

MLP-based SCM (70%)  $f$  : Activation  $\circ$  Linear



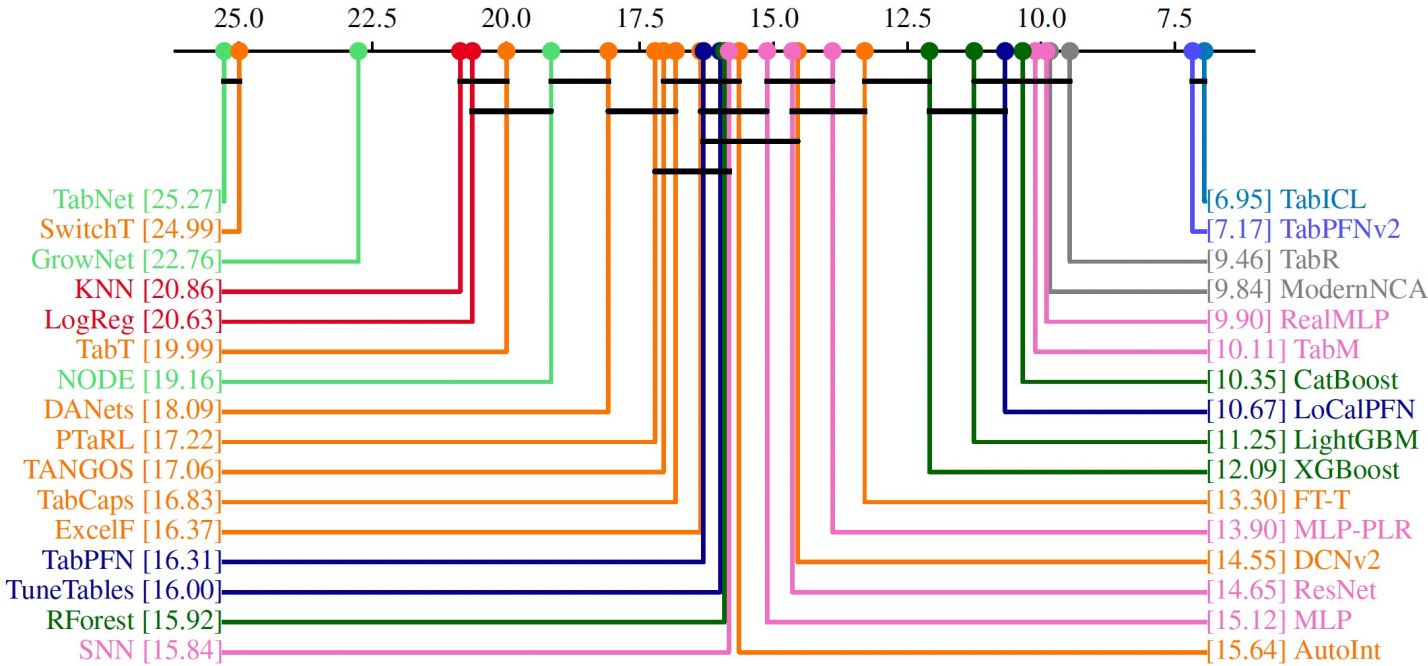
The weights and biases of the linear layer are randomly sampled.

Tree-based SCM (30%)  $f$  : XGBoost [Chen et al. 2016]



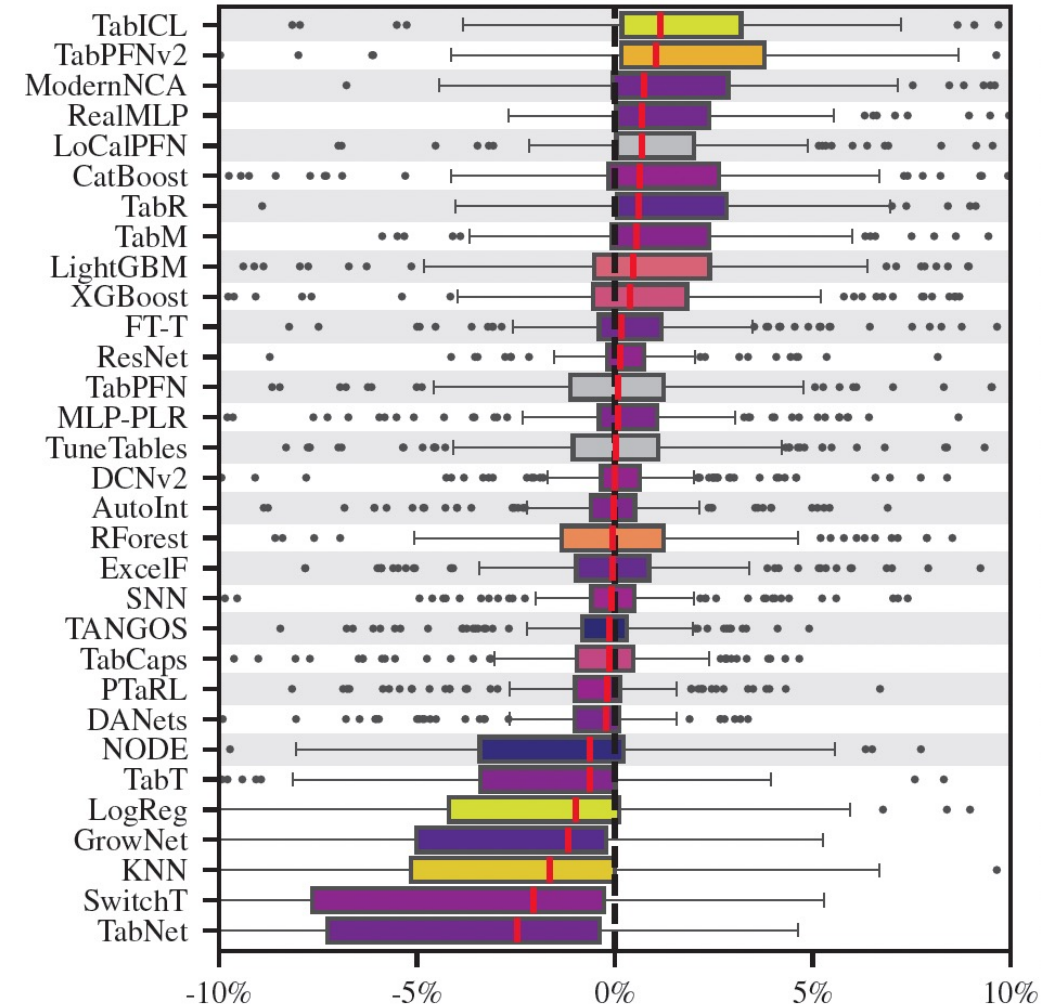
XGBoost is trained on fake targets drawn from Gaussian noise.

Aspect		TabPFNv2	TabICL
Architecture	Attention	Alternating column / row attentions	Column -> Row -> ICL
	Collapse Issue	Feature grouping and random feature vectors	RoPE
	Label Fusion	Early (input layer)	Late (for ICL)
Pre-training	SCM	Growing network with redirection (Code not open-source)	Layered structure
	# of datasets	82 million	130 million
	Curriculum learning	✗ # of samples $\leq 2,048$	✓ # of samples 1K -> 60K
Scalability	# of samples	$\leq 10K$	100K (only 5GB GPU)
	# of classes	$\leq 10$	Any (Hierarchical classification)

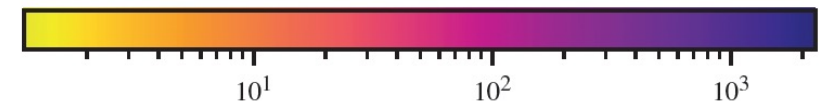


Average accuracy ranking over 171 classification datasets  
( $\leq 10$  classes) from the TALENT benchmark

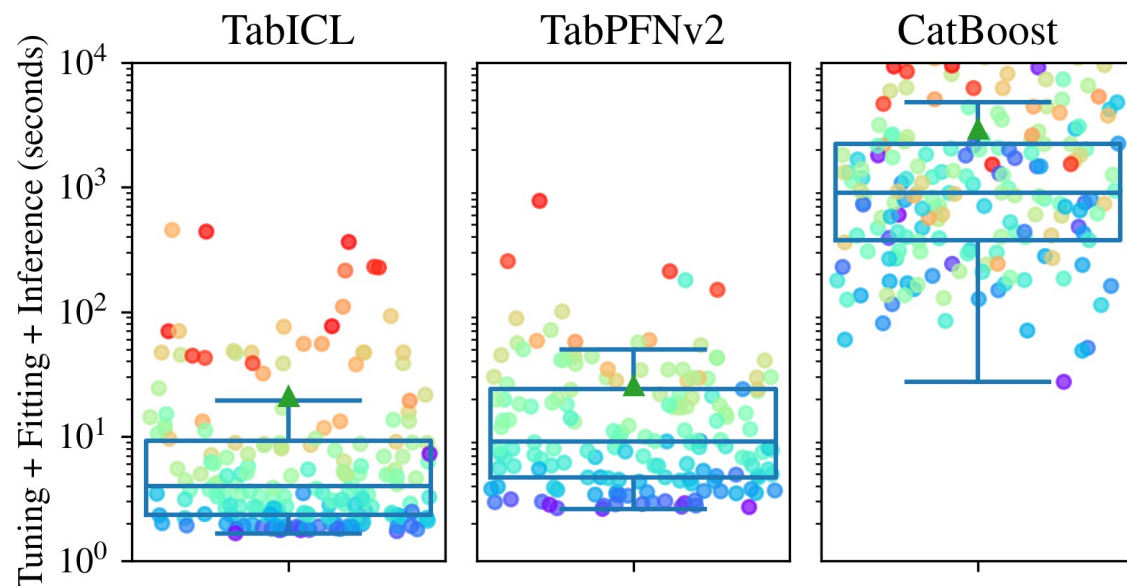
[Ye et al. 2024]



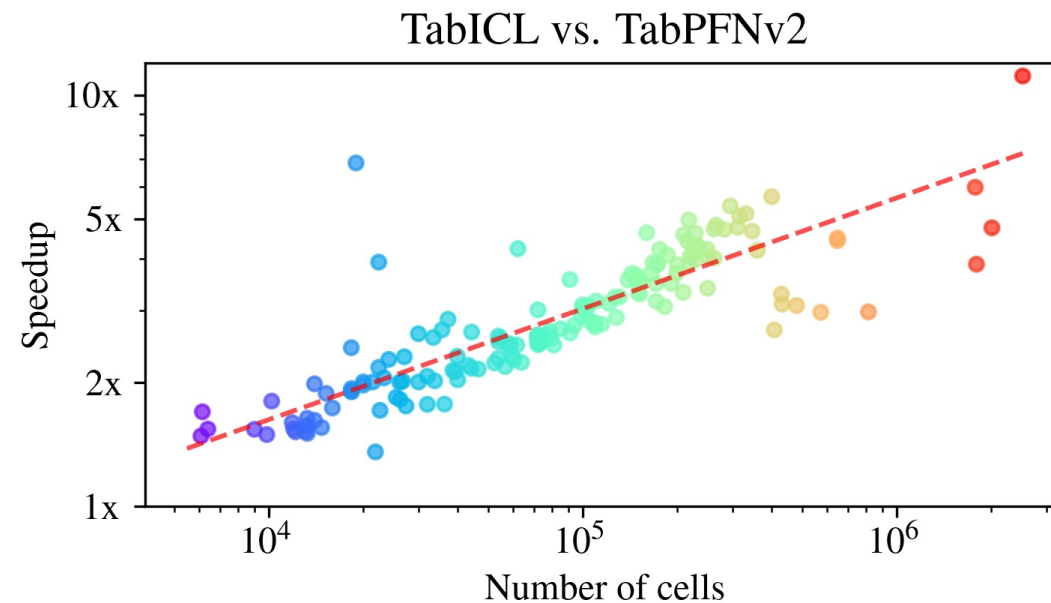
Relative accuracy improvement over MLP ( $\uparrow$ )



Geom. mean train time [s] per 1K samples

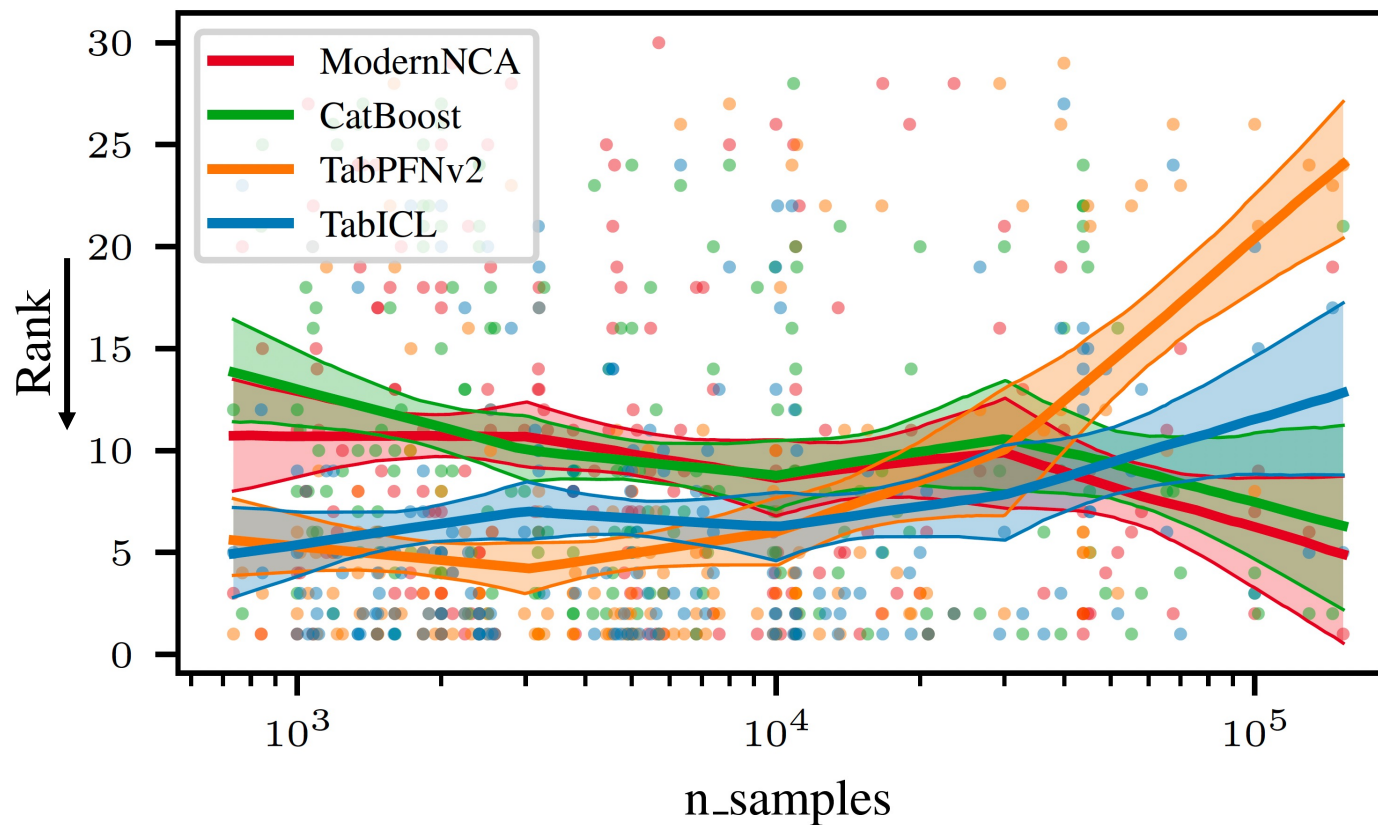


Time comparison per dataset

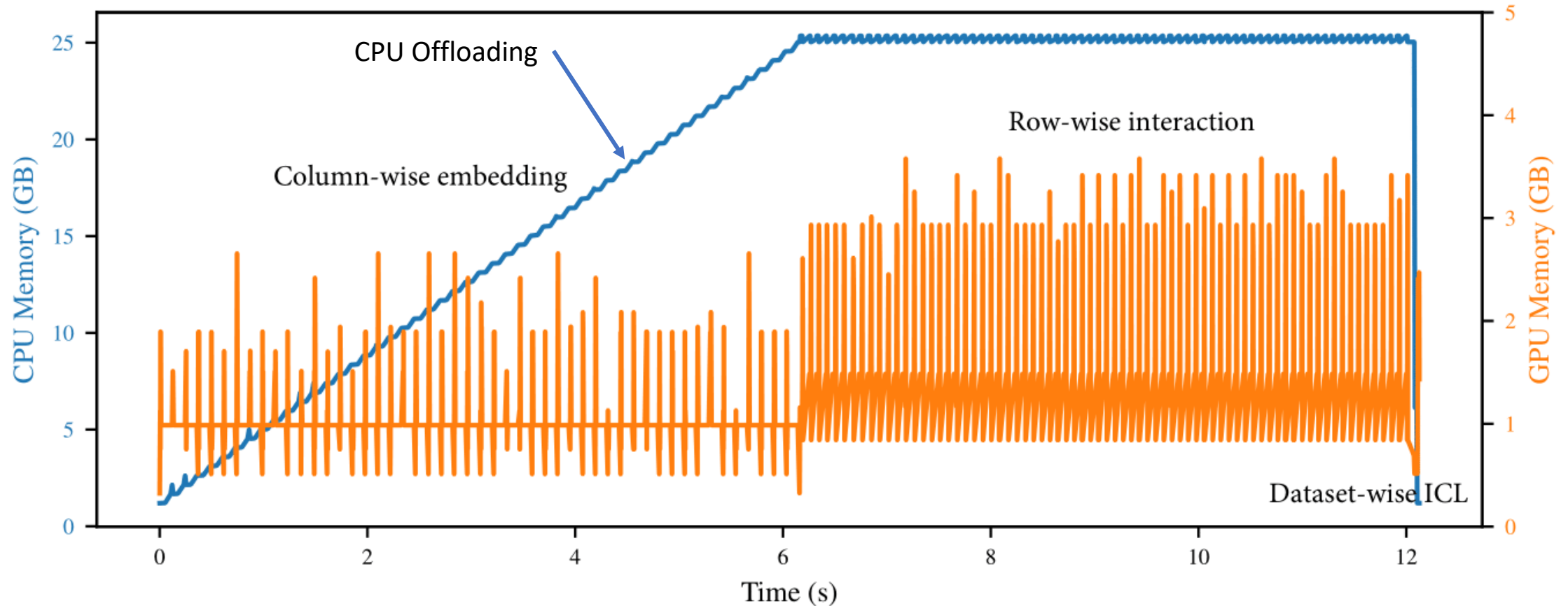


Speedup of TabICL vs. TabPFNv2





Model rankings vs. sample size



CPU and GPU memory during inference for a dataset with 100K samples and 500 features (< **5GB GPU memory**)

- Tabular foundation models can be scaled to an order of magnitude larger data!
  - ❖ Expressive yet computationally efficient architecture
  - ❖ Large-scale pre-training via curriculum learning by gradually increasing data size
  - ❖ Memory-efficient inference through CPU offloading, adaptive batching, etc.
  
- Faster and better than TabPFNv2 despite
  - ❖ Simpler preprocessing
  - ❖ Simpler prior (no random graphs, no categorical vectors/discretization)
  
- Open-source everything : <https://github.com/soda-inria/tabicl>

*Thanks for your attention !*