

TabPFGen - Tabular Data Generation with TabPFN

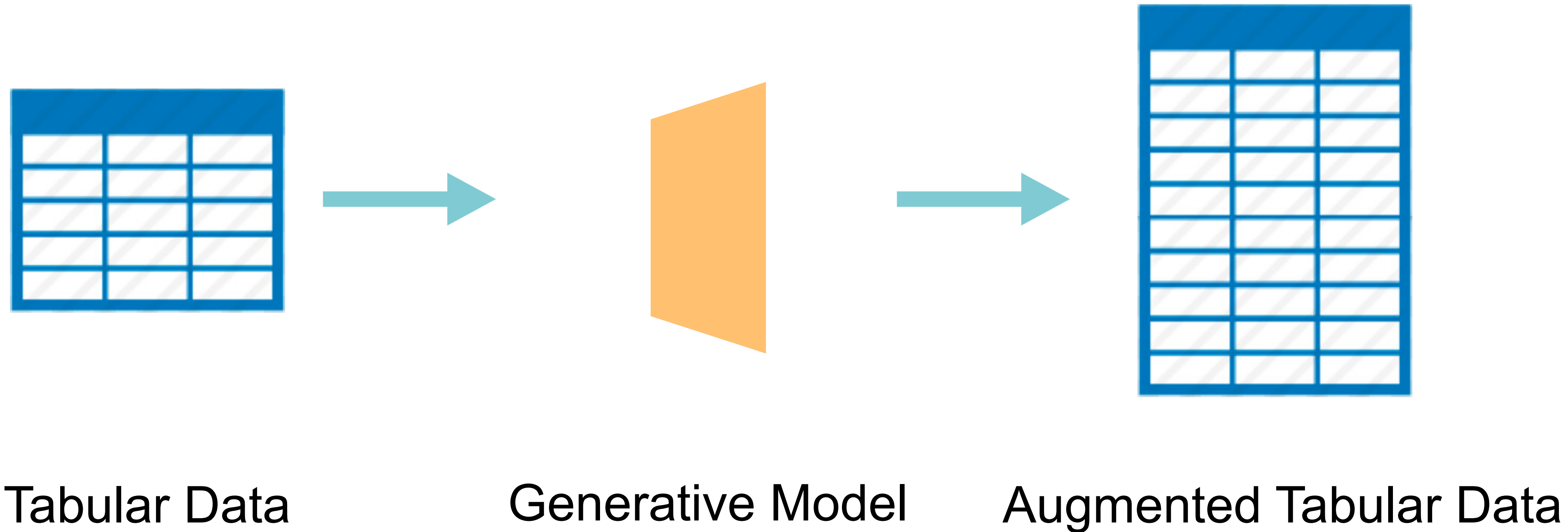
Junwei Ma, Apoorv Dankar, George Stein, Guangwei Yu, Anthony Caterini

15 Dec, 2023

layer6

Tabular Data Generation

- The task is to generate more data according to the underlying data generation process



Goals of Tabular Data Generation

- **Augmentation** - for performance
- **Replacement** - for privacy
- **Class Balancing** - for imbalanced data
- **Imputation** - for missing data
- **Data Summarization** - for redundant data

Challenges

- **Generalization across datasets**

- Past works usually involve unstable and time-consuming training and hyper-parameter tuning
- The problem exacerbates when the data size is small

- **Inductive bias**

- The inductive bias in tabular data is not clear especially with small size data
- This is in contrast to images and text where the inductive bias is commonly exploited

Motivation - TabPFN

- TabPFN (Hollmann et al.) outperforms other tree-based and deep learning methods according to independent study (McElfresh et al.)
 - Can we harness this high-performing discriminative model for generation?
- TabPFN is exposed to a large number of data generation processes and inductive biases
 - Can we utilize the discriminative power of TabPFN for generation?

Harnessing TabPFN - Energy-based Models

- Energy-based models (EBMs) background
 - EBMs parameterize a density using its unnormalized log-density function

- $p_{\theta}(x) = \frac{\exp(-E_{\theta}(x))}{Z_{\theta}}$, where $Z_{\theta} = \int_x \exp(-E_{\theta}(x)) dx$

- E_{θ} can be any function or network
- How do we define the energy function for generation?

Harnessing TabPFN - the Energy Function

- Grathwohl et al. proposed to use classifier outputs as the energy function

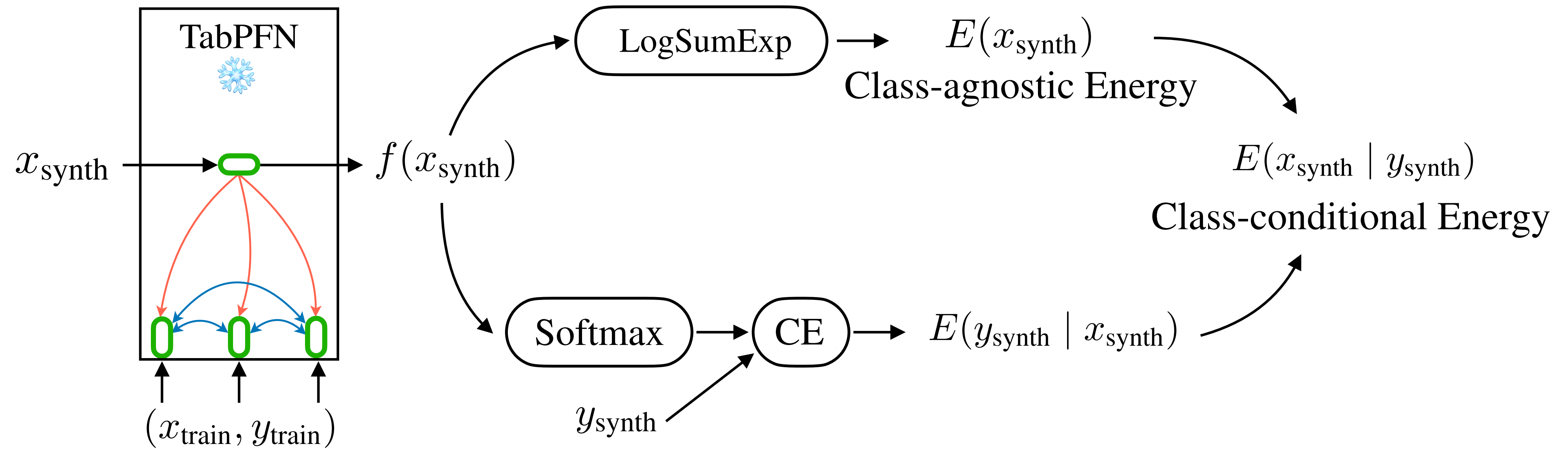
- $$E_{\theta}(x) = -\log \sum_y \exp(f_{\theta}(x)[y]),$$

- where $f_{\theta}(x)$ is the classifier logit and $[\]$ indicates indexing function

- To sample from this EBM, we can use Stochastic Gradient Langevin Dynamics (SGLD)

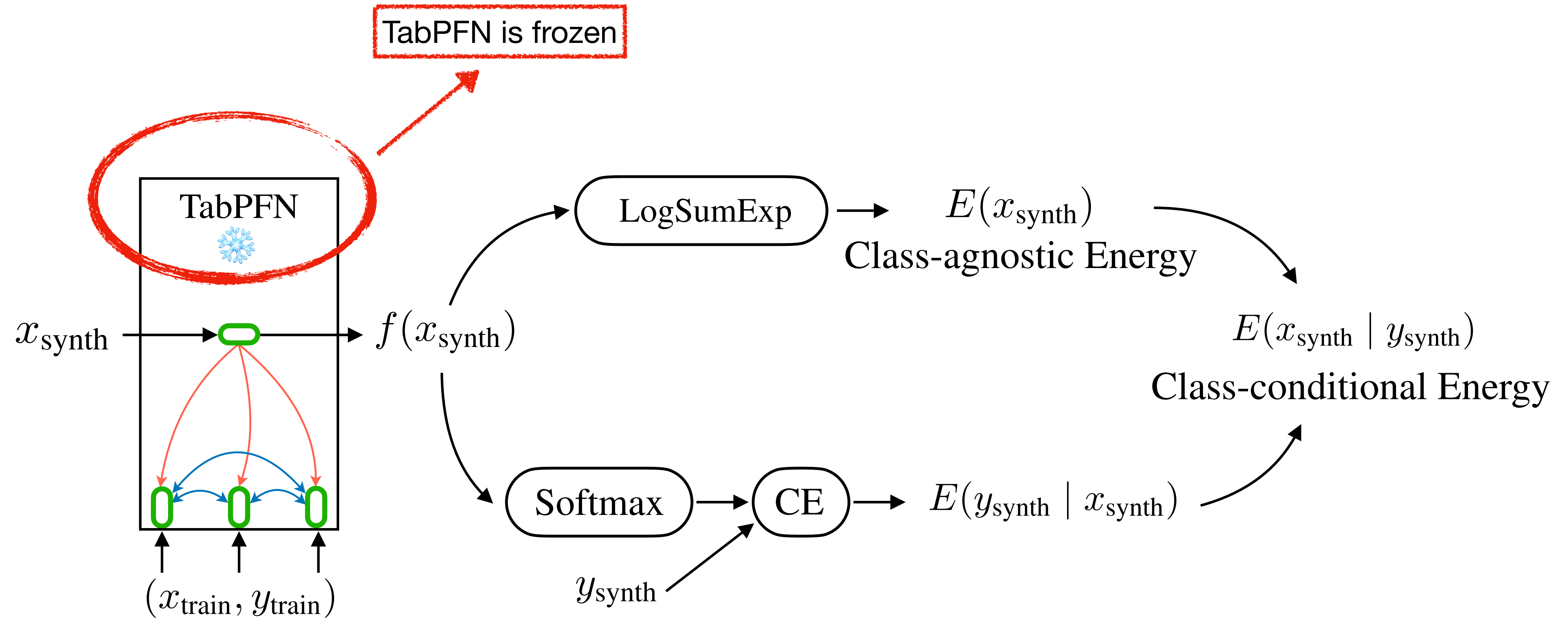
layer 6

TabPFGen



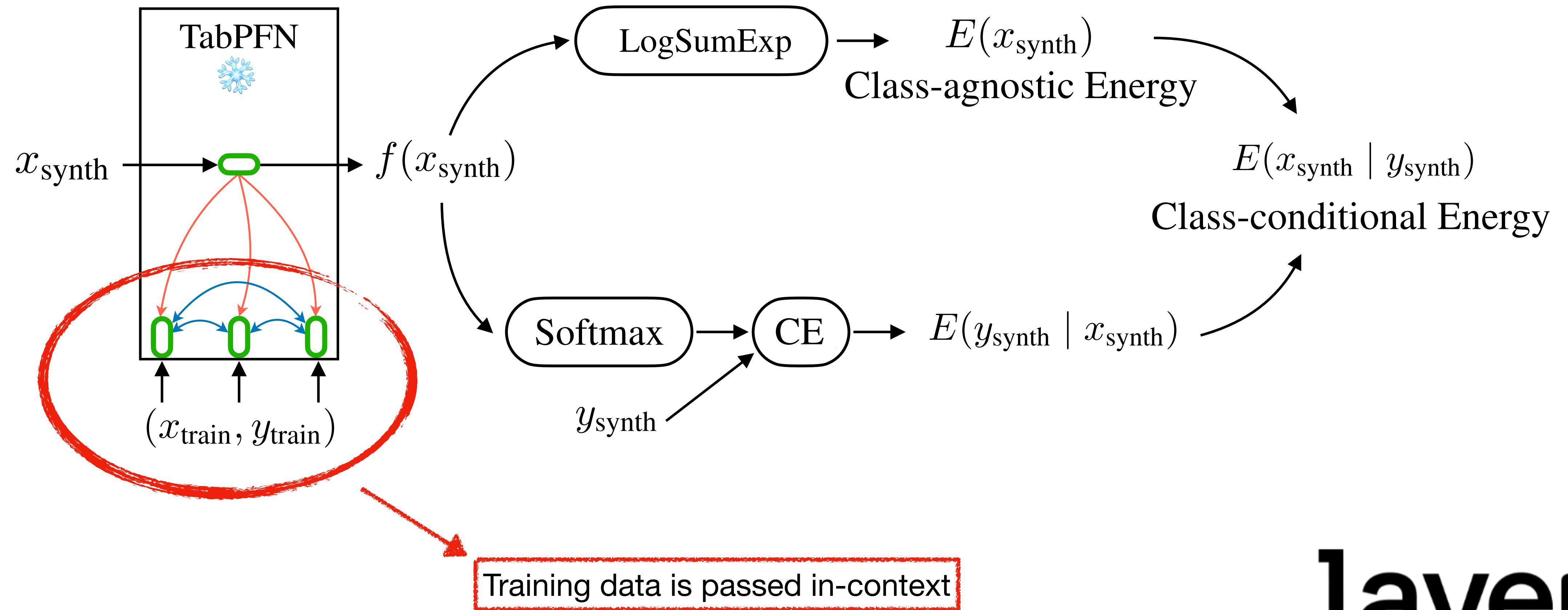
layer 6

TabPFN



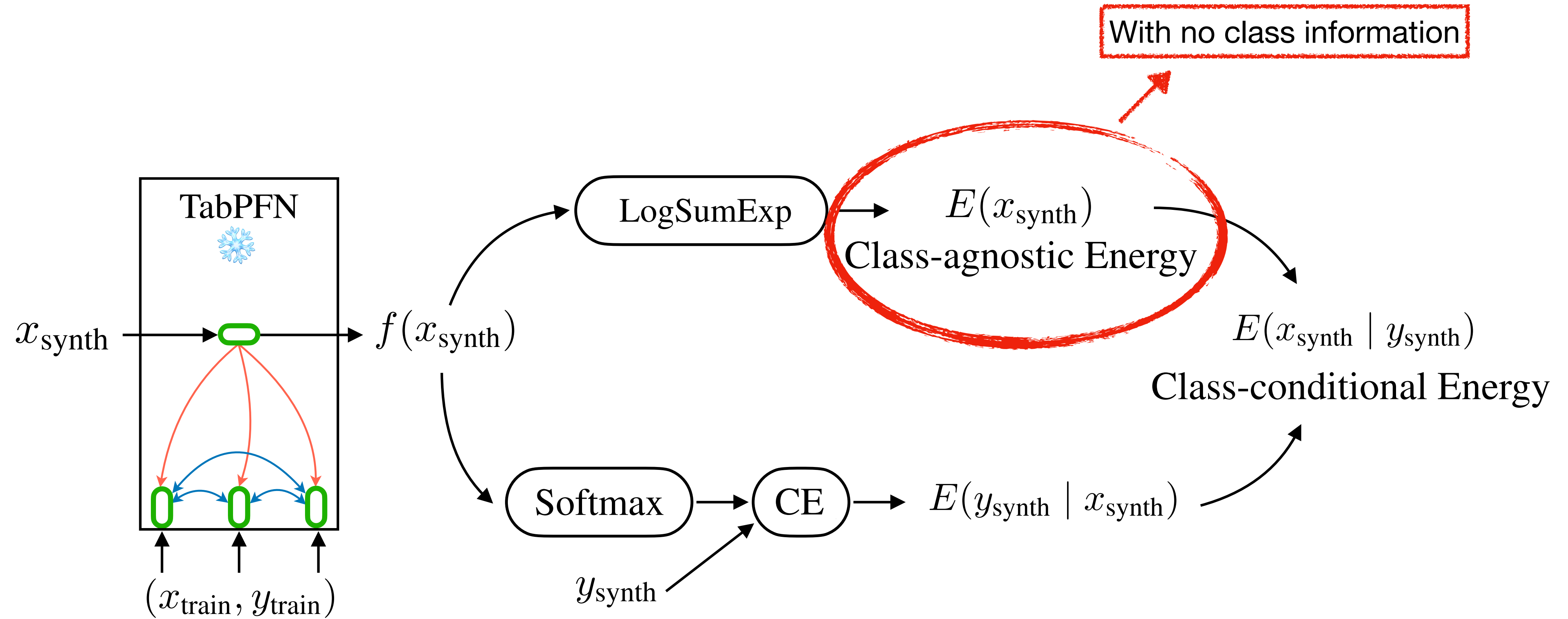
layer 6

TabPFGen



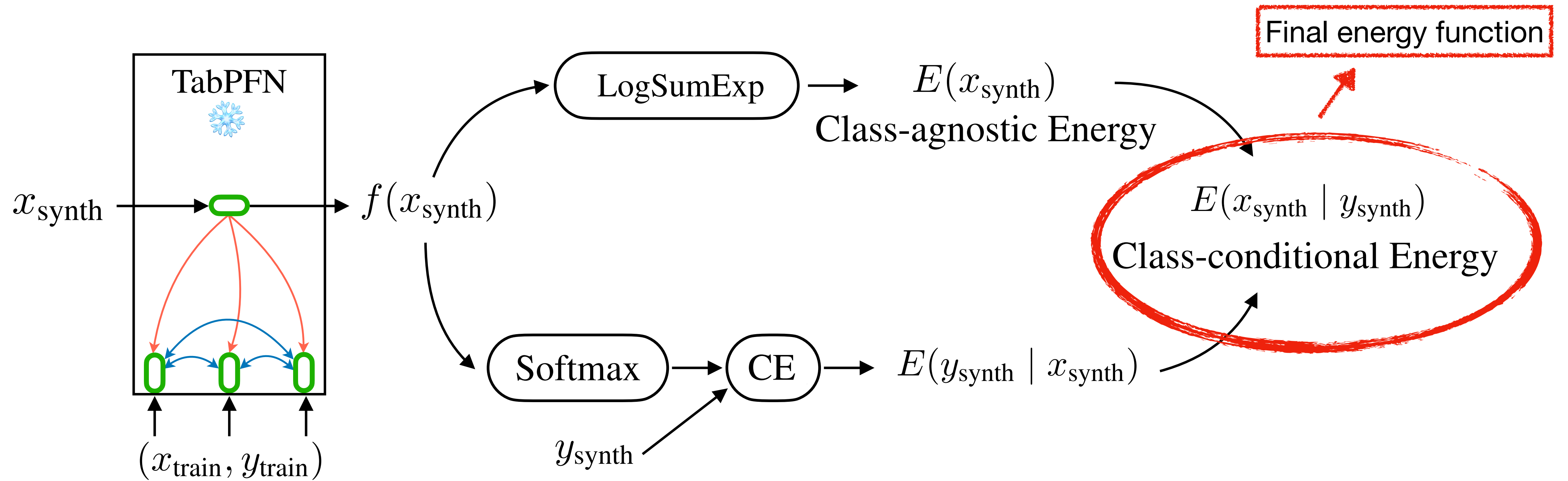
layer 6

TabPFGen



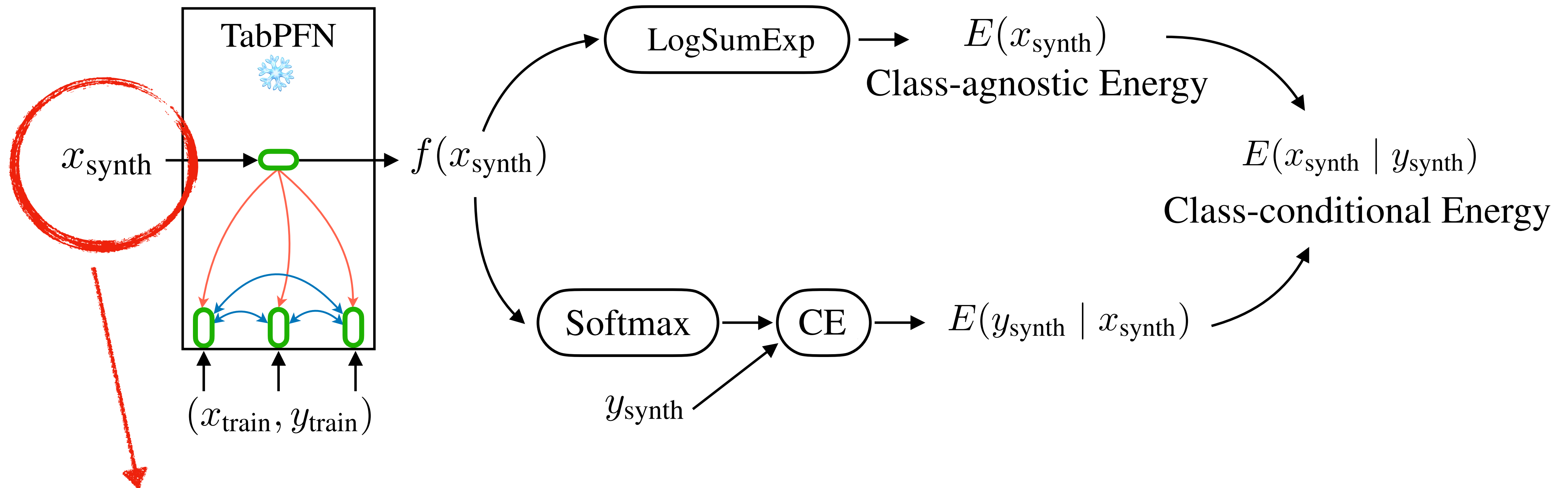
layer 6

TabPFGen



layer 6

TabPFGen

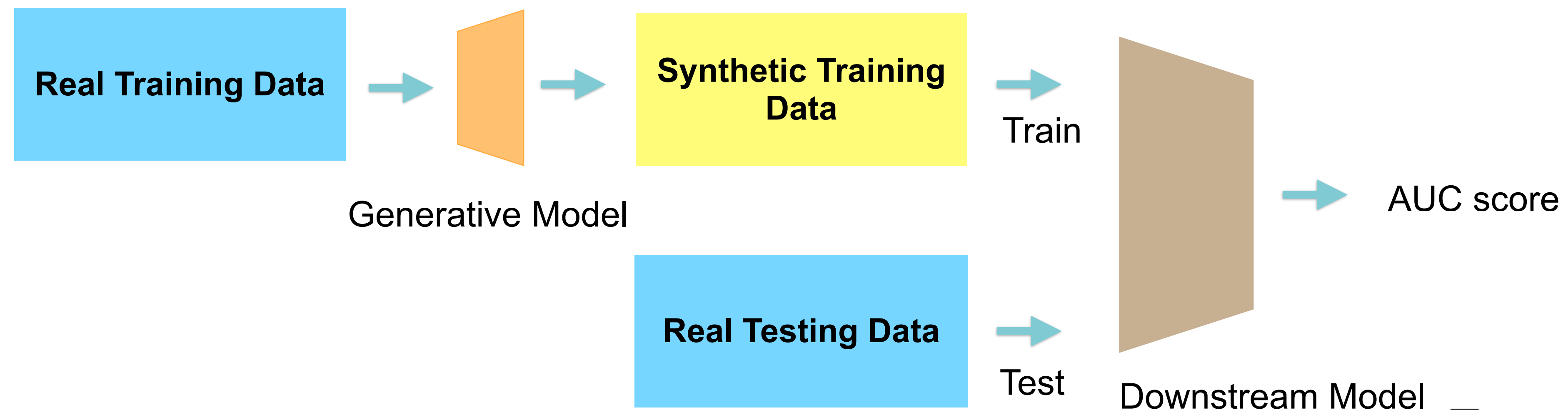


In-context learning generation with SGLD:
- Frozen TabPFN's discriminative power is inverted for generation
- No additional training or hyper-parameter tuning required

layer 6

Experimental Setup

- Use synthetic data to augment, replace or class balance
- Test set is used to evaluate downstream models performance with synthetic data
- Downstream models include: XGB, RF, LR and TabPFN



Results - Augmentation and Replacement

Augmentation

Model	Original	SMOTE	CTGAN	TVAE	NF	RTVAE	TabDDPM	TabPFGen
XGB	$0.924 \pm 3e-4$	$0.926 \pm 3e-4$	$0.912 \pm 2e-4$	$0.914 \pm 7e-4$	$0.912 \pm 4e-4$	$0.917 \pm 3e-4$	$0.927 \pm 3e-4$	$0.934 \pm 3e-4$
RF	$0.906 \pm 3e-4$	$0.906 \pm 2e-3$	$0.898 \pm 1e-3$	$0.904 \pm 1e-3$	$0.894 \pm 3e-4$	$0.907 \pm 2e-3$	$0.911 \pm 7e-4$	$0.912 \pm 4e-4$
LR	$0.920 \pm 7e-4$	$0.914 \pm 3e-3$	$0.904 \pm 3e-3$	$0.909 \pm 6e-3$	$0.901 \pm 9e-4$	$0.906 \pm 8e-3$	$0.885 \pm 3e-4$	$0.921 \pm 2e-4$
TabPFN	$0.934 \pm 2e-3$	$0.927 \pm 1e-3$	$0.930 \pm 1e-3$	$0.931 \pm 1e-3$	$0.928 \pm 3e-4$	$0.932 \pm 1e-3$	$0.929 \pm 5e-4$	$0.935 \pm 3e-4$
XGB	N/A	$0.907 \pm 4e-4$	$0.842 \pm 8e-4$	$0.858 \pm 2e-3$	$0.700 \pm 6e-4$	$0.795 \pm 9e-4$	$0.812 \pm 3e-4$	$0.927 \pm 3e-4$
RF	N/A	$0.894 \pm 1e-3$	$0.837 \pm 6e-4$	$0.844 \pm 5e-4$	$0.676 \pm 2e-3$	$0.774 \pm 3e-4$	$0.814 \pm 9e-4$	$0.906 \pm 6e-4$
LR	N/A	$0.893 \pm 2e-3$	$0.843 \pm 6e-4$	$0.873 \pm 1e-3$	$0.722 \pm 3e-3$	$0.854 \pm 7e-4$	$0.876 \pm 3e-4$	$0.920 \pm 1e-3$
TabPFN	N/A	$0.920 \pm 8e-4$	$0.888 \pm 4e-4$	$0.887 \pm 3e-4$	$0.705 \pm 2e-3$	$0.862 \pm 1e-3$	$0.894 \pm 7e-4$	$0.934 \pm 2e-4$

layer 6

Results - Augmentation and Replacement

Model	Original	SMOTE	CTGAN	TVAE	NF	RTVAE	TabDDPM	TabPFGen
XGB	$0.924 \pm 3e-4$	$0.926 \pm 3e-4$	$0.912 \pm 2e-4$	$0.914 \pm 7e-4$	$0.912 \pm 4e-4$	$0.917 \pm 3e-4$	$0.927 \pm 3e-4$	$0.934 \pm 3e-4$
RF	$0.906 \pm 3e-4$	$0.906 \pm 2e-3$	$0.898 \pm 1e-3$	$0.904 \pm 1e-3$	$0.894 \pm 3e-4$	$0.907 \pm 2e-3$	$0.911 \pm 7e-4$	$0.912 \pm 4e-4$
LR	$0.920 \pm 7e-4$	$0.914 \pm 3e-3$	$0.904 \pm 3e-3$	$0.909 \pm 6e-3$	$0.901 \pm 9e-4$	$0.906 \pm 8e-3$	$0.885 \pm 3e-4$	$0.921 \pm 2e-4$
TabPFN	$0.934 \pm 2e-3$	$0.927 \pm 1e-3$	$0.930 \pm 1e-3$	$0.931 \pm 1e-3$	$0.928 \pm 3e-4$	$0.932 \pm 1e-3$	$0.929 \pm 5e-4$	$0.935 \pm 3e-4$
XGB	N/A	$0.907 \pm 4e-4$	$0.842 \pm 8e-4$	$0.858 \pm 2e-3$	$0.700 \pm 6e-4$	$0.795 \pm 9e-4$	$0.812 \pm 3e-4$	$0.927 \pm 3e-4$
RF	N/A	$0.894 \pm 1e-3$	$0.837 \pm 6e-4$	$0.844 \pm 5e-4$	$0.676 \pm 2e-3$	$0.774 \pm 3e-4$	$0.814 \pm 9e-4$	$0.906 \pm 6e-4$
LR	N/A	$0.893 \pm 2e-3$	$0.843 \pm 6e-4$	$0.873 \pm 1e-3$	$0.722 \pm 3e-3$	$0.854 \pm 7e-4$	$0.876 \pm 3e-4$	$0.920 \pm 1e-3$
TabPFN	N/A	$0.920 \pm 8e-4$	$0.888 \pm 4e-4$	$0.887 \pm 3e-4$	$0.705 \pm 2e-3$	$0.862 \pm 1e-3$	$0.894 \pm 7e-4$	$0.934 \pm 2e-4$

Replacement

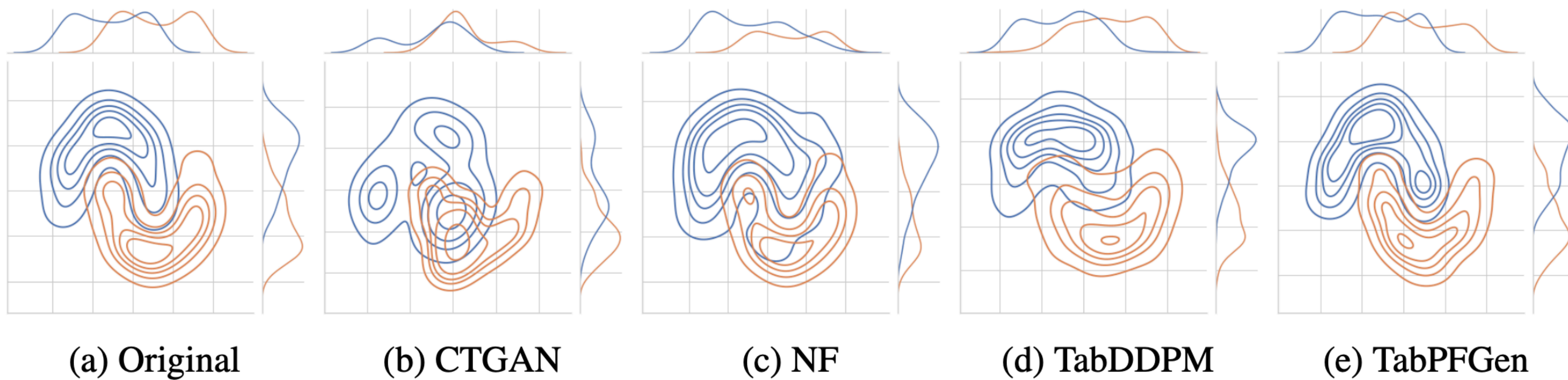
layer 6

Results - Class Balancing

Dataset	Original	Sampling	SMOTE	CTGAN	TVAE	NF	RTVAE	TabDDPM	TabPFGen
KC	0.823	0.875	0.872	0.859	0.848	0.862	0.866	0.805	0.877
PC	0.824	0.811	0.836	0.827	0.835	0.825	0.841	0.825	0.841
BL	0.731	0.757	0.756	0.743	0.714	0.755	0.706	0.774	0.767
CL	0.925	0.935	0.949	0.793	0.771	0.795	0.909	0.915	0.955
DI	0.837	0.832	0.832	0.831	0.813	0.832	0.837	0.843	0.844

layer 6

Results - Qualitative



layer 6

A black water bottle with a white cap and a black carabiner is placed on a weathered wooden post on a beach. The bottle has the word "layer6" printed vertically in white. The background shows a sunset over the ocean with a person visible in the distance.

Thank you!

jeremy@layer6.ai