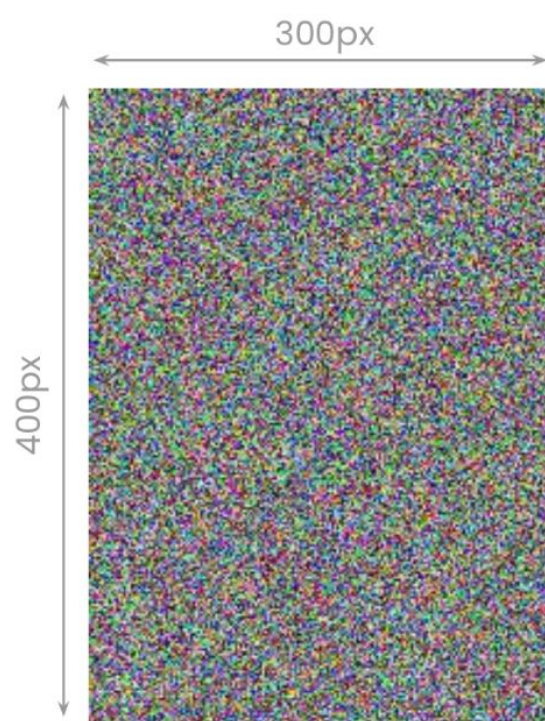


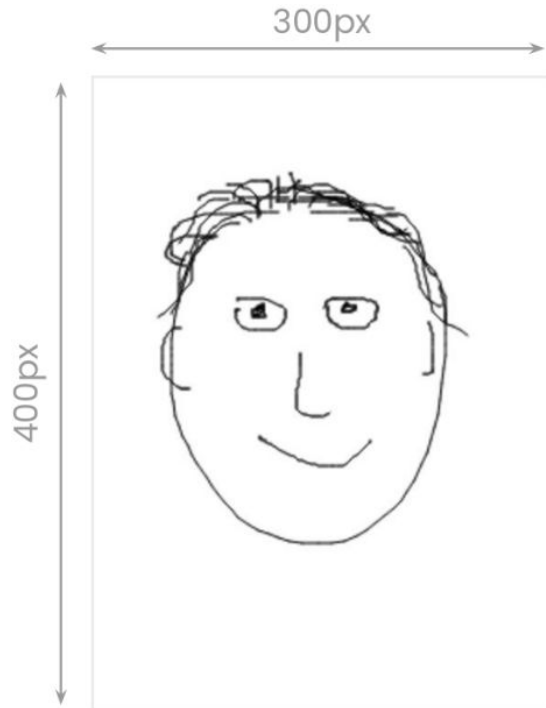
TabularARGN

A Framework for High-Quality Synthetic Data Generation

What is Synthetic Data



random data



self-generated data



model-generated data
rule-based

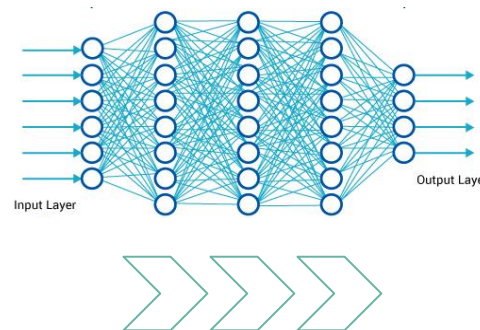


AI-generated data
"data-based"

What is Tabular Synthetic Data

NAME	AGE	GENDER	ITEM	EUR	DATE	TIME
Mary	25y	female	Book	12€	4/2/19	8:12
John	72y	male	Pizza	34€	4/2/19	18:12
...						
Bill	18y	male	Swim	6€	4/4/19	10:02
Bill	18y	male	Shoes	123€	4/4/19	12:32

Real Data



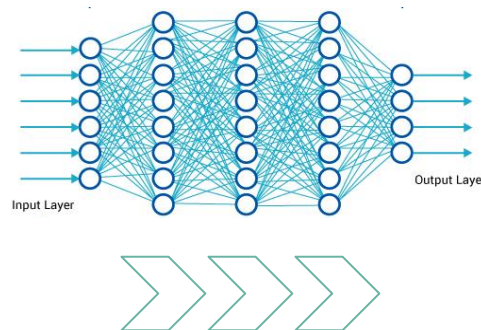
NAME	AGE	GENDER	ITEM	EUR	DATE	TIME
Kim	29y	female	Amazon	236€	4/4/19	12:32
Kim	29y	female	Zalando	36€	4/4/19	18:58
...						
Brian	82y	male	Beer	6€	4/2/19	21:32
Sue	24y	female	Sushi	12€	4/2/19	21:32

Synthetic Data

Privacy-preserving Tabular Synthetic Data

NAME	AGE	GENDER	ITEM	EUR	DATE	TIME
Mary	25y	female	Book	12€	4/2/19	8:12
John	72y	male	Pizza	34€	4/2/19	18:12
...						
Bill	18y	male	Swim	6€	4/4/19	10:02
Bill	18y	male	Shoes	123€	4/4/19	12:32

Real Data



Private

NAME	AGE	GENDER	ITEM	EUR	DATE	TIME
Kim	29y	female	Amazon	236€	4/4/19	12:32
Kim	29y	female	Zalando	36€	4/4/19	18:58
...						
Brian	82y	male	Beer	6€	4/2/19	21:32
Sue	24y	female	Sushi	12€	4/2/19	21:32

Synthetic Data

Open

Tabular ARGN - Auto-Regressive Generative Networks

arXiv > cs > arXiv:2501.12012

Search...

Help

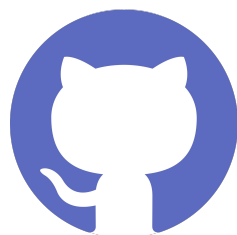
Computer Science > Machine Learning

[Submitted on 21 Jan 2025 (v1), last revised 6 Feb 2025 (this version, v2)]

TabularARGN: A Flexible and Efficient Auto-Regressive Framework for Generating High-Fidelity Synthetic Data

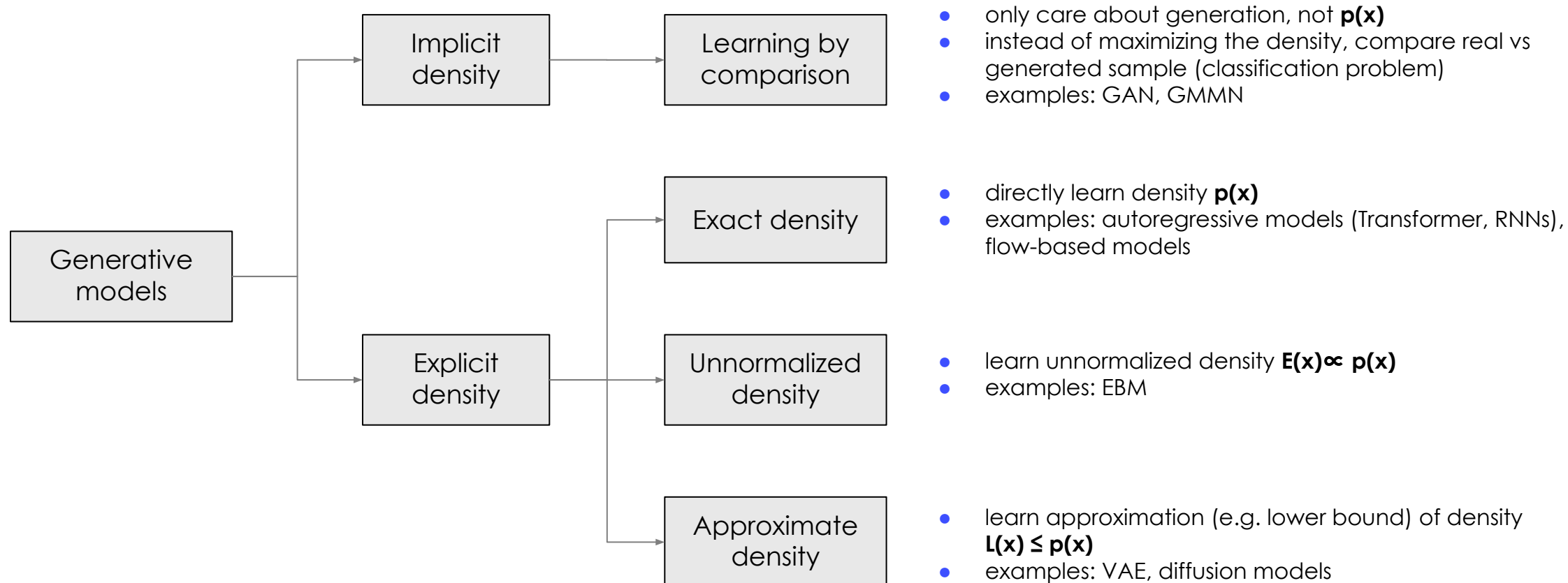
Paul Tiwald, Ivona Krchova, Andrey Sidorenko, Mariana Vargas Vieyra, Mario Scriminaci, Michael Platzer

Synthetic data generation for tabular datasets must balance fidelity, efficiency, and versatility to meet the demands of real-world applications. We introduce the Tabular Auto-Regressive Generative Network (TabularARGN), a flexible framework designed to handle mixed-type, multivariate, and sequential datasets. By training on all possible conditional probabilities, TabularARGN supports advanced features such as fairness-aware generation, imputation, and conditional generation on any subset of columns. The framework achieves state-of-the-art synthetic data quality while significantly reducing training and inference times, making it ideal for large-scale datasets with diverse structures. Evaluated across established benchmarks, including realistic datasets with complex relationships, TabularARGN demonstrates its capability to synthesize high-quality data efficiently. By unifying flexibility and performance, this framework paves the way for practical synthetic data generation across industries.



Tabular ARGN is implemented in the [Synthetic Data SDK](#)

Taxonomy of deep generative models



Flat Model

Fixed (column) Order Training Phase

patients data set:

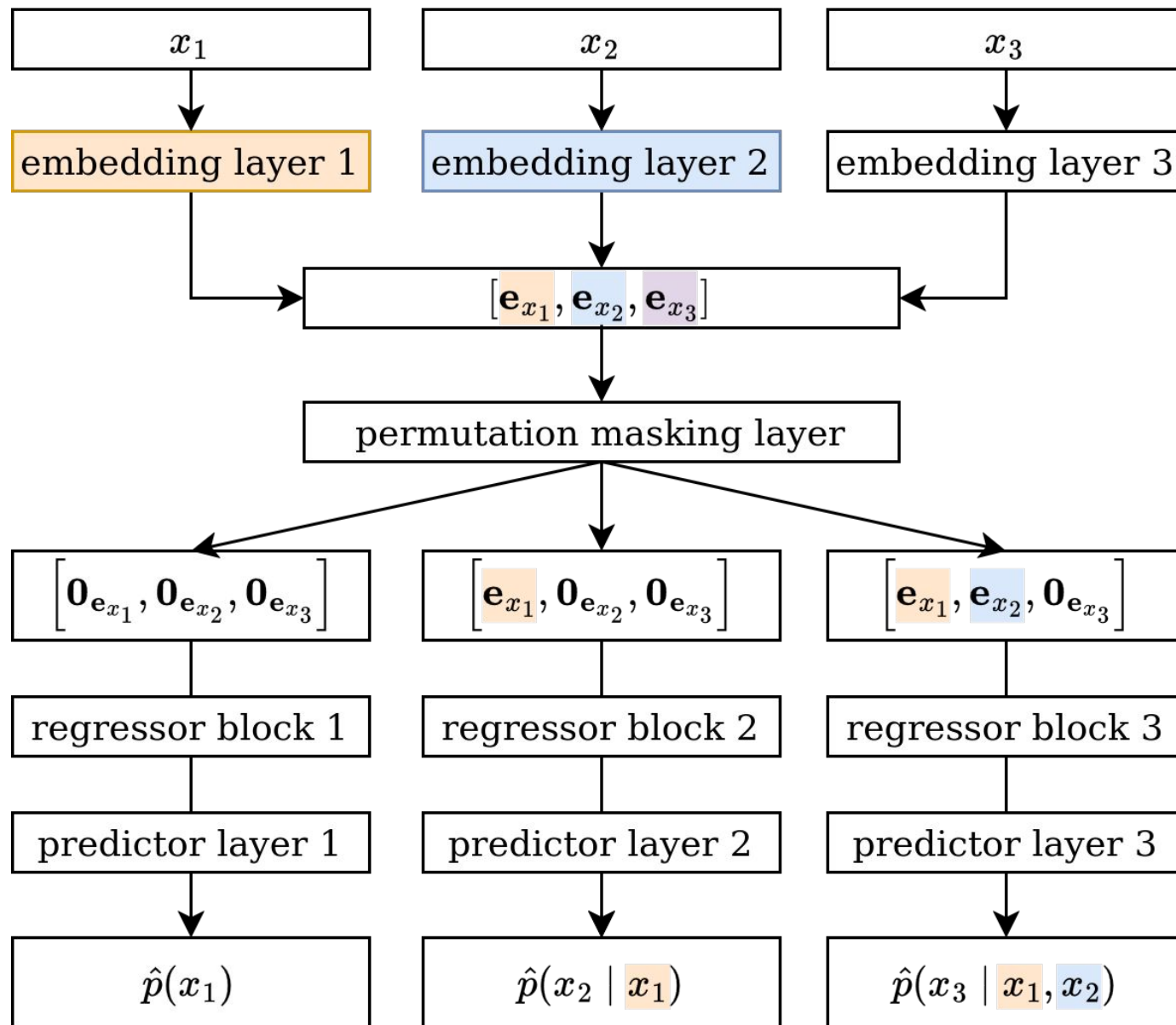
x_1 - age

x_2 - gender

x_3 - blood type

loss function:

$$\max_{\theta} \sum_{i=1}^D \log p_{\theta}(x_i \mid x_{<i})$$



Flat Model

Any (column) Order Training Phase

patients data set:

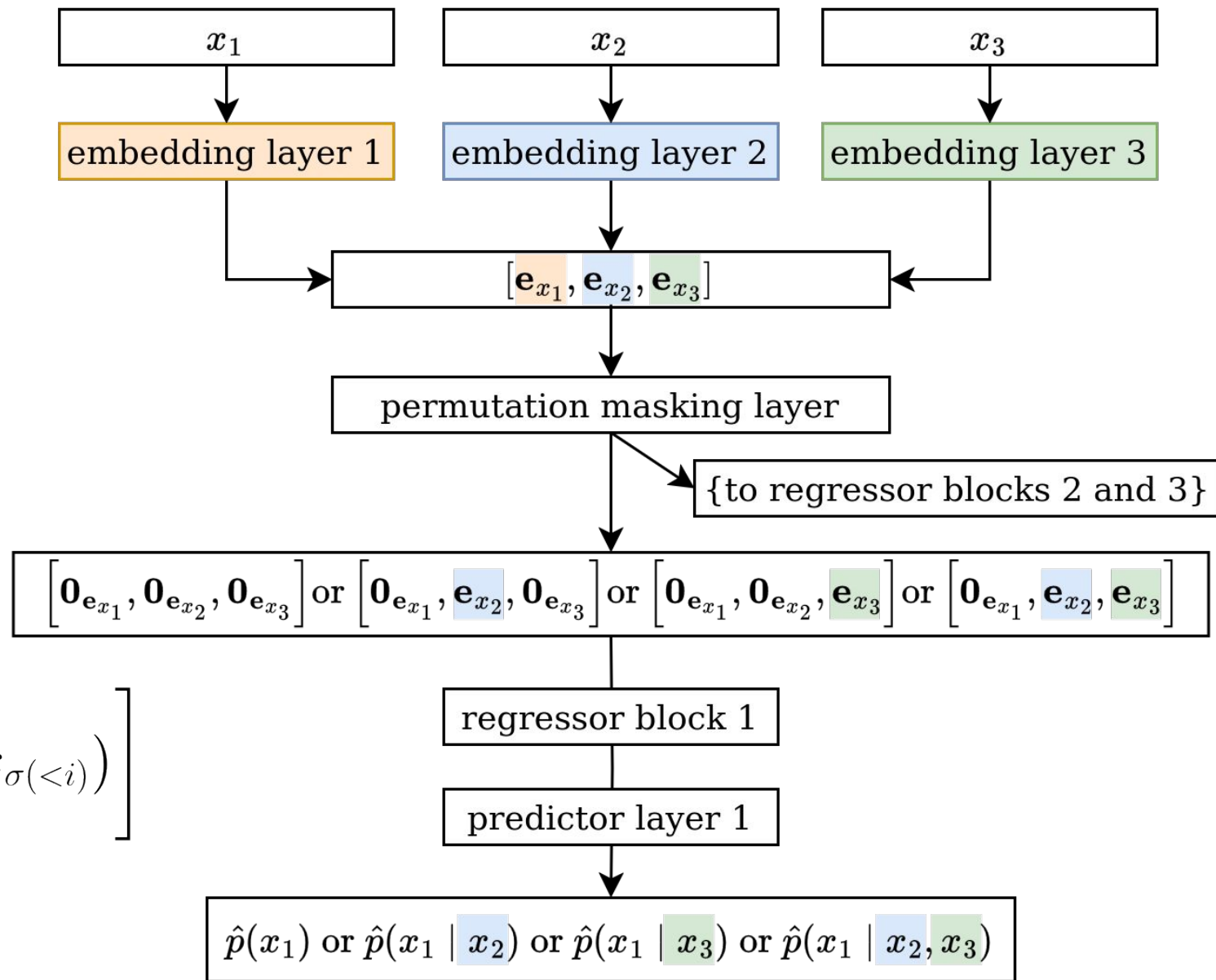
x_1 - age

x_2 - gender

x_3 - blood type

loss function:

$$\max_{\theta} \mathbb{E}_{\sigma \in \text{Uniform}(S_D)} \left[\sum_{i=1}^D \log p_{\theta}(x_{\sigma(i)} \mid x_{\sigma(<i)}) \right]$$



Flat Model

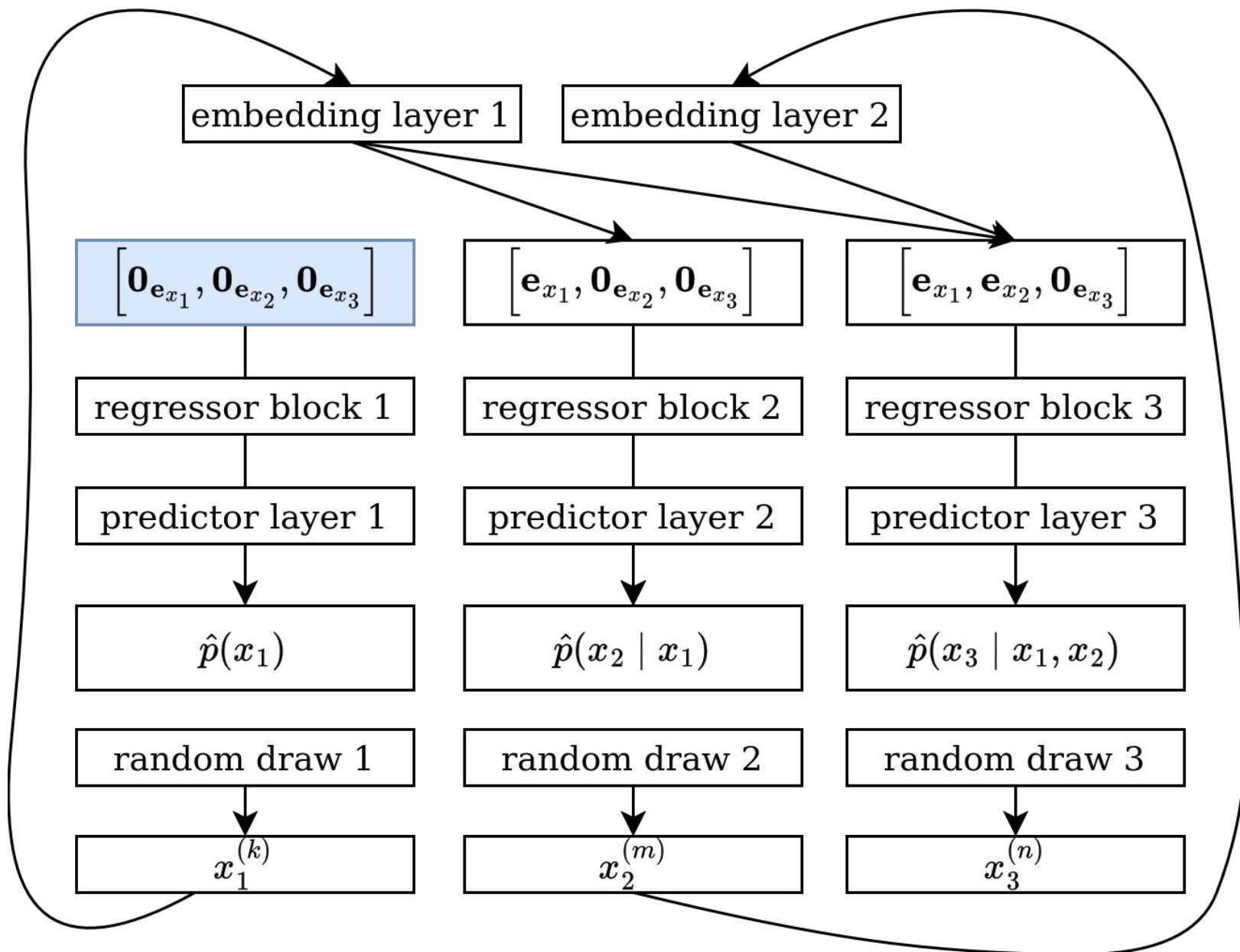
Generation Phase

patients data set:

x_1 - age

x_2 - gender

x_3 - blood type

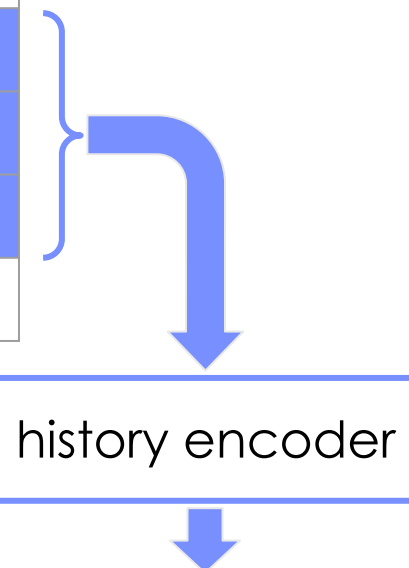


Sequential Model - doctor visits

auto-regressive along the column and the time dimensions

← time

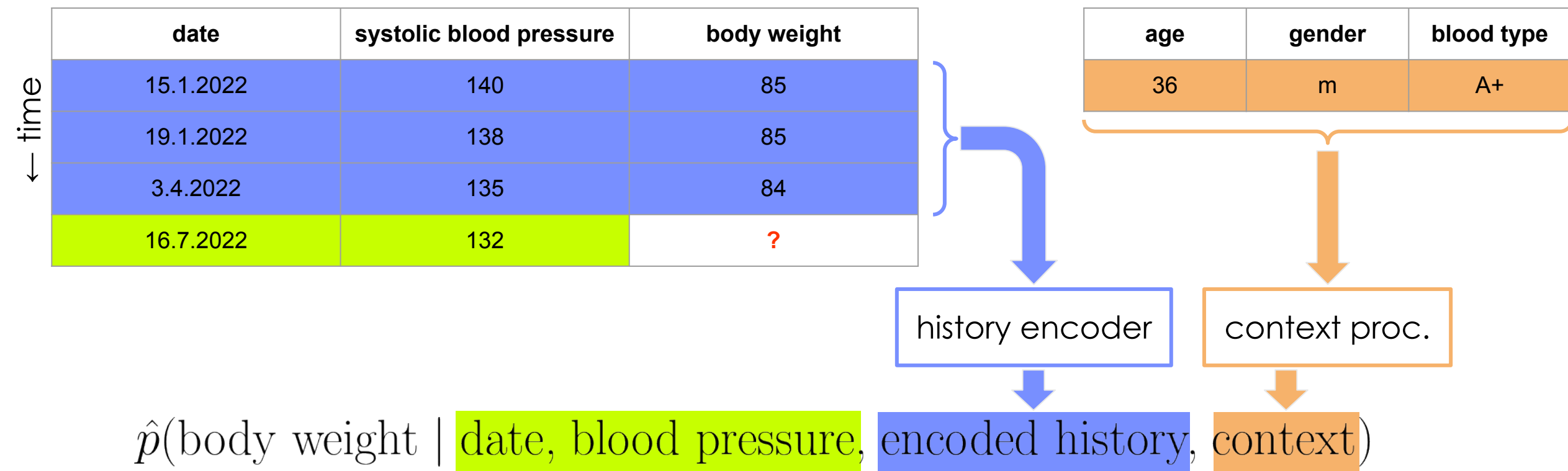
date	systolic blood pressure	body weight
15.1.2022	140	85
19.1.2022	138	85
3.4.2022	135	84
16.7.2022	132	?



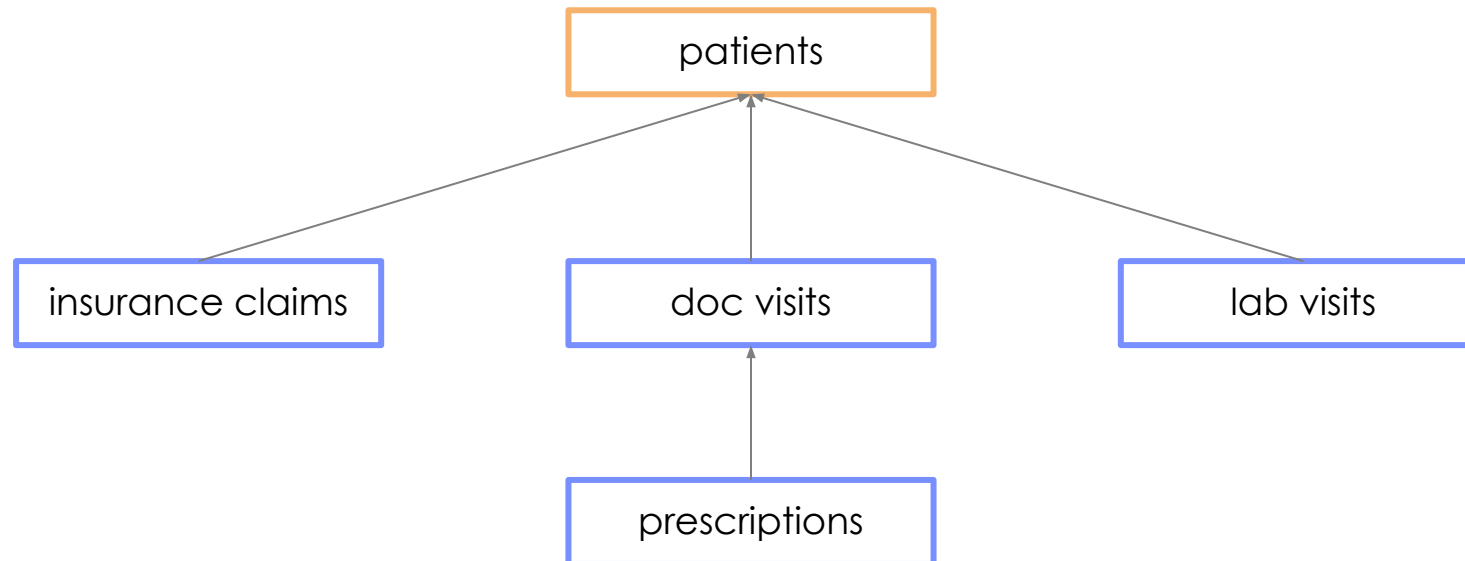
$$\hat{p}(\text{body weight} \mid \text{date, blood pressure, encoded history})$$

Sequential Model with context

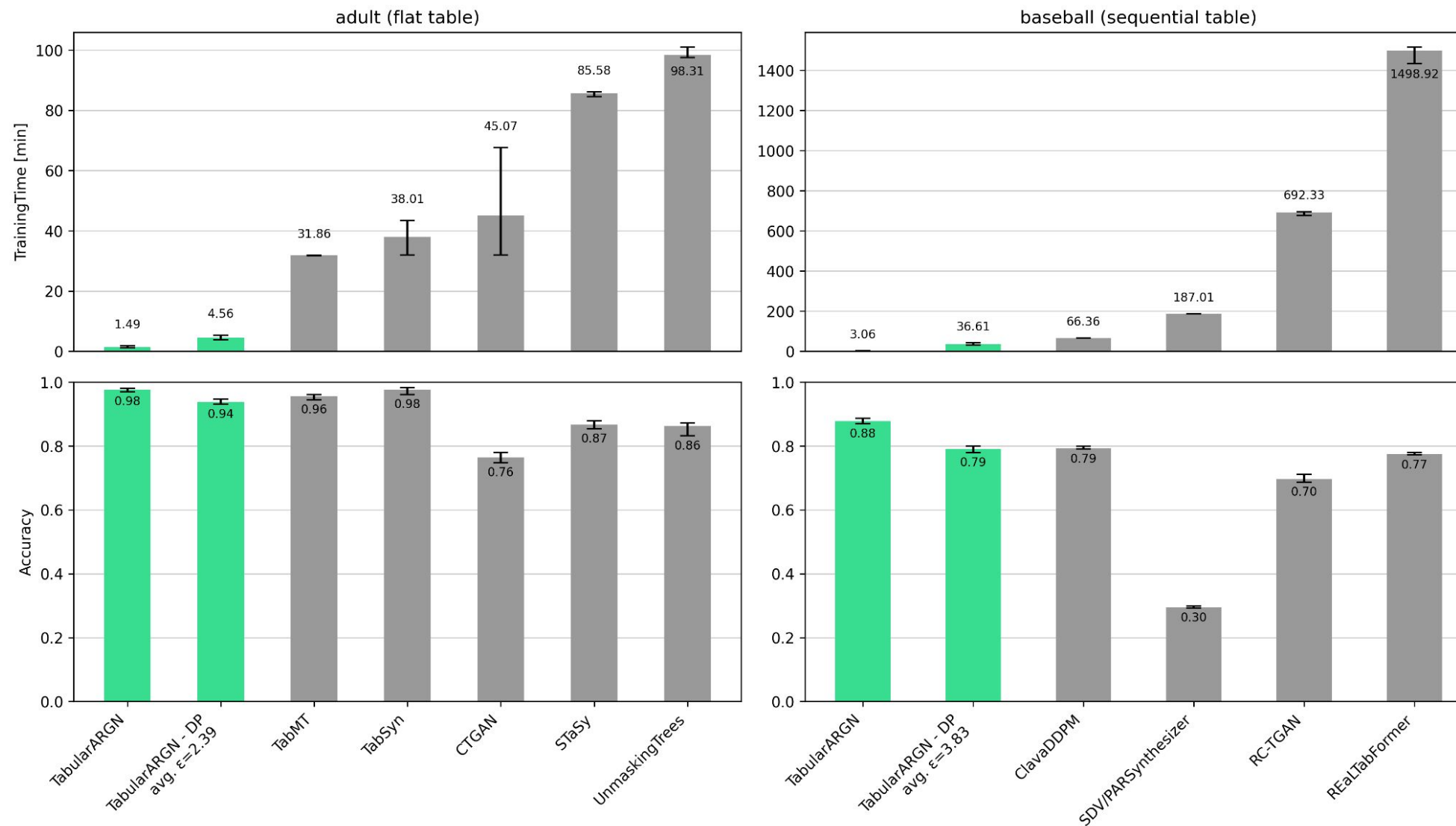
auto-regressive along the column, time, and table dimensions



Flexible context allows for synthesis of multi-table setups



Benchmarking TabularARGN against SOTA methods



Text-extension allows for leveraging the power of LLMs

combining TabularARGN with text models

date	systolic blood pressure	body weight	diagnosis
15.1.2022	140	85	"The patient presents with ..."

TabularARGN

LanguageARGN

training data:

```
{ "input": "15.1.2022; 140; 85", "output": "The patient presents with ..." }
```

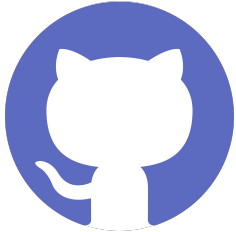

We are hiring

Developer Advocate (Remote)

Permanent employee, Full-time · EMEA Remote

We are looking for a captivating Developer Advocate to help us spread the word, engage with the Python / Data / AI community, and create some fun content. Somewhere at the **intersection of a Data Enthusiast, Open Source Advocate, and a Content Creator.**

Tabular ARGN - Auto-Regressiv Generative Networks



Tabular ARGN is implemented in the [Synthetic Data SDK](#)



for more infos, check out our blog

[Synthetic Behavioral Data](#)

[Synthetic Geo Data](#)

[Differentially Private Synthetic Data](#)

[Fair Synthetic Data](#)

[Synthetic Data Benchmarks](#)

[JRC Report on Synthetic Data](#)

[AI-based Re-Identification Attacks](#)

[Privacy Assessment of Synthetic Data](#)

and many more