


Multi-Modal Tabular Deep Learning with PyTorch Frame

Akihiro Nitta

akihiro@kumo.ai

 /pyg-team/pytorch-frame



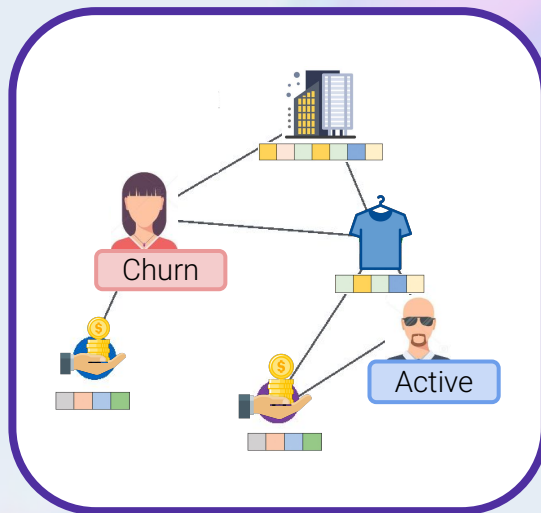
```
pip install pytorch-frame
```

Table of Contents

- Challenges in Modern Tabular Learning
- PyTorch Frame – Multi-modal tabular deep learning
- Single table to multiple tables

Tree-based models are dominant – with limitations

- GBDTs are focused on numerical and categorical features
 - Modern tabular data have text and images
- Integrating GBDT with deep learning models is non-trivial
 - training tabular models with image and text encoders
 - training recommendation models with GNNs end-to-end





PyTorch Frame

PyTorch Frame

A modular framework for tabular learning

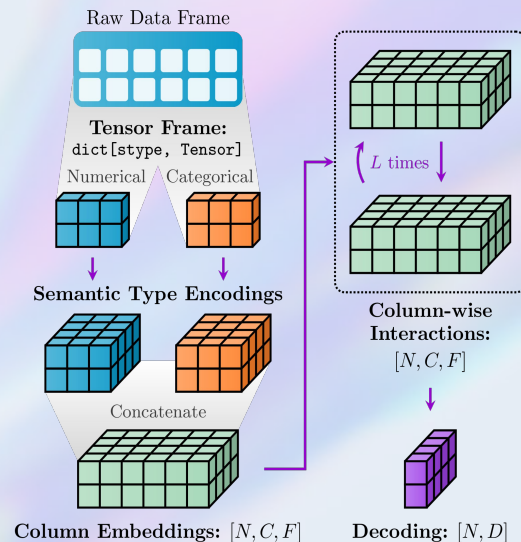
- ✓ Multi-modal features
numerical, categorical, multi-categorical,
image, text, embedding, timestamp
- ✓ Modular design covering various existing models
e.g., Trompt, ExcelFormer, FT-Transformer
- ✓ Integration with foundation models & GNNs



 Best Paper Award at TRL@NeurIPS 2024 

```
pip install pytorch-frame
```

 pyg-team/pytorch-frame



Data Materialization

During data materialization, 🏠 PyTorch Frame:

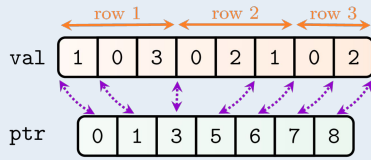
- **encodes a raw data frame** into a TensorFrame
- **computes statistics**, e.g., mean, standard deviation, count of category elements

TensorFrame is a tensor-based data structure

- stores columns of the same type in a tensor
- stores sparse features efficiently, e.g., via MultiNestedTensor

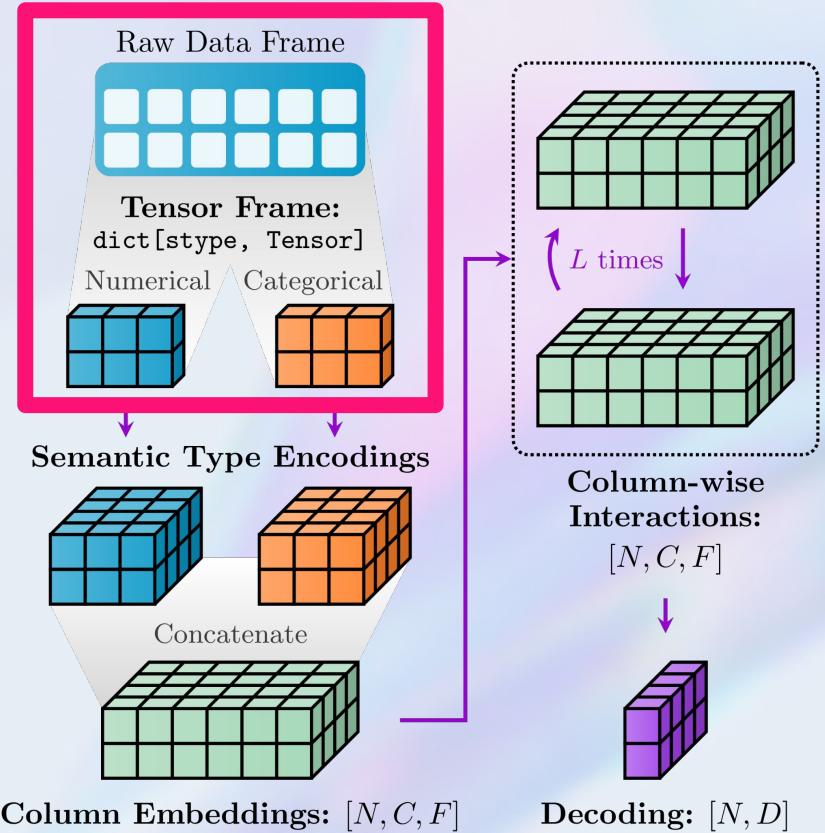
| | |
|-----|-----|
| 1 | 0,3 |
| 0,2 | 1 |
| 0 | 2 |

Raw Data Frame



Compressed Representation

Data Materialization



Data Materialization and Semantic Types

📖 PyTorch Frame materializes columns of different *semantic types* as follows:

- **numerical** – Pass through of floating-point columns
- **categorical** – Map categories to indices
- **multicategorical** – Map multiple categories into indices of varying length
- **timestamp** – Map timestamps to integers of year, month, day, hour, minute, ...
- **embedding** – Pass through of pre-computed embeddings
- **text_tokenized** – Tokenize text into a list of integers of varying length
- **text_embedded** – Pre-compute text vectors via external text models
- **image_embedded** – Pre-compute image vectors via external image models

Native Integration with Foundation Models



OpenAI



cohere



Hugging Face

VOYAGE AI

Data Materialization and Semantic Types

📖 PyTorch Frame materializes columns of different *semantic types* as follows:

- **numerical** – Pass through of floating-point columns
- **categorical** – Map categories to indices
- **multicategorical** – Map multiple categories into indices of varying length
- **timestamp** – Map timestamps to integers of year, month, day, hour, minute, ...
- **embedding** – Pass through of pre-computed embeddings
- **text_tokenized** – Tokenize text into a list of integers of varying length
- **text_embedded** – Pre-compute text vectors via external text models
- **image_embedded** – Pre-compute image vectors via external image models

Native Integration with Foundation Models



Hugging Face

VOYAGE AI

Data Materialization and Text Columns

🔥 PyTorch Frame provides two ways to handle text columns:

- Pre-encode text during **materialization** once to avoid computing the same embeddings

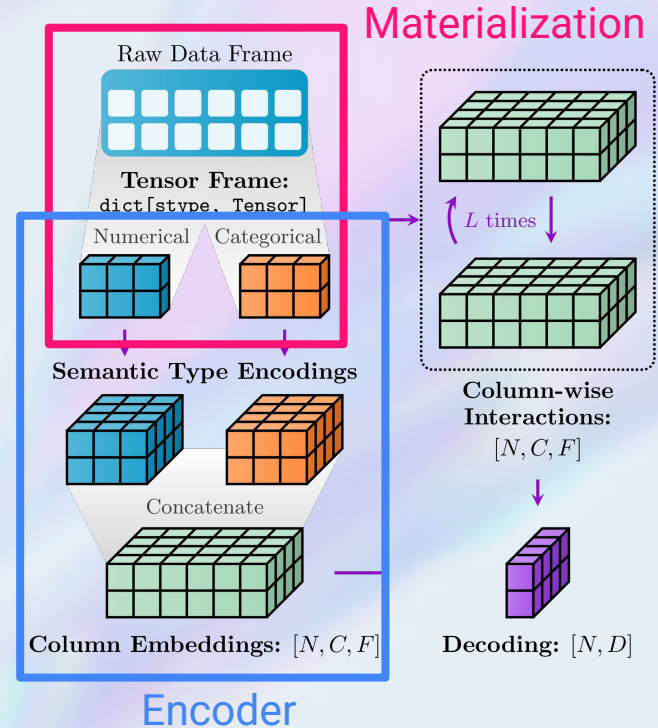


“text_embedded”

- Tokenize text during **materialization** once
- Finetune embeddings in **encoder**



“text_tokenized”



Data Materialization and Semantic Types

📖 PyTorch Frame materializes columns of different *semantic types* as follows:

- **numerical** – Pass through of floating-point columns
- **categorical** – Map categories to indices
- **multicategorical** – Map multiple categories into indices of **varying length**
- **timestamp** – Map timestamps to integers of year, month, day, hour, minute, ...
- **embedding** – Pass through of pre-computed embeddings
- **text_tokenized** – Tokenize text into a list of integers of **varying length**
- **text_embedded** – Pre-compute text vectors via external text models
- **image_embedded** – Pre-compute image vectors via external image models

Native Integration with Foundation Models



OpenAI



cohere



Hugging Face

VOYAGE AI

Sparse Features in

PyTorch Frame stores variable-length features efficiently via `torch_frame.data.MultiNestedTensor`

Size of dense representation:

$$N * \text{max_entries} * \text{sizeof}(\text{dtype})$$

Size of sparse representation:

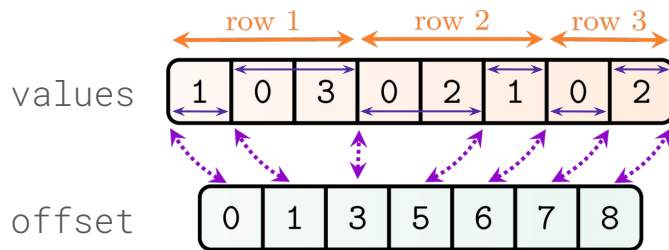
$$N * \text{max_entries} * \text{sizeof}(\text{dtype}) * \text{density} \\ + (N + 1) * \text{sizeof}(\text{int64})$$

Example:

- dtype is int64, N is 1,000, max_entries is 1,000, density is 10%
- dense: 8,000,000 bytes
- sparse: 8,08,008 bytes (only 10% of dense version!)

| | col 1 | col 2 |
|-------|-------|-------|
| row 1 | 1 | 0,3 |
| row 2 | 0,2 | 1 |
| row 3 | 0 | 2 |

Raw Data Frame



Compressed Representation

Three-stage Model Architecture

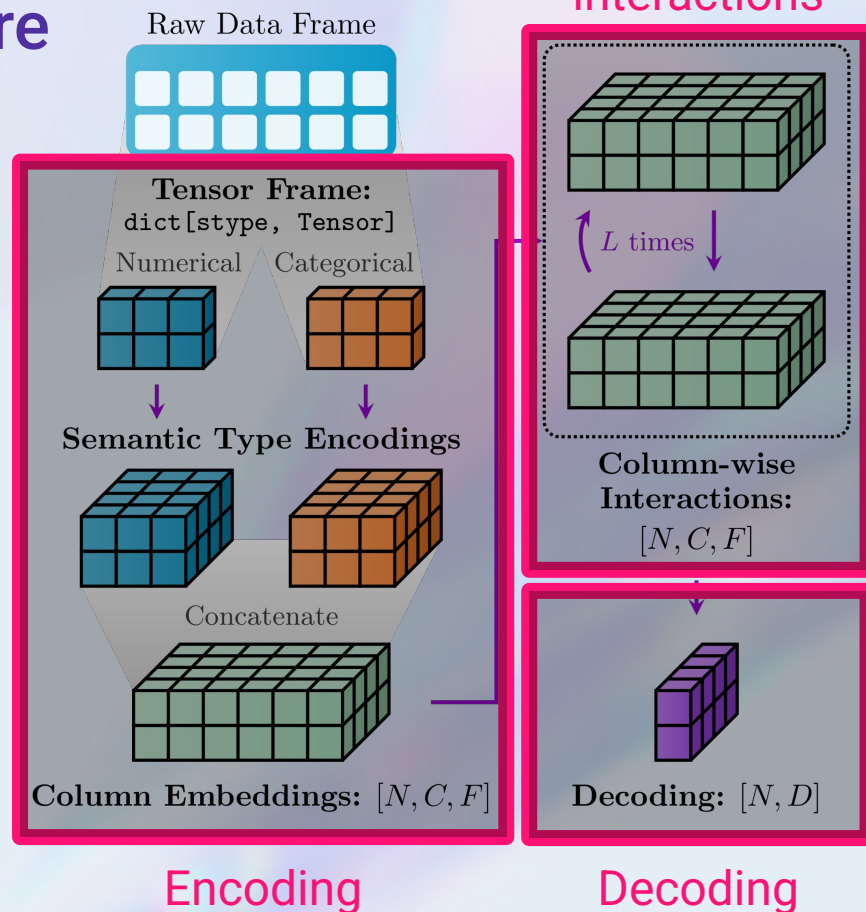
Encoding embeds each column independently

- missing data imputation
- normalization
- embedding lookup
- cyclic encoding
- positional encoding
- ...

Column-wise interaction performs message passing across columns.

Decoding summarizes column embeddings to obtain row embeddings.

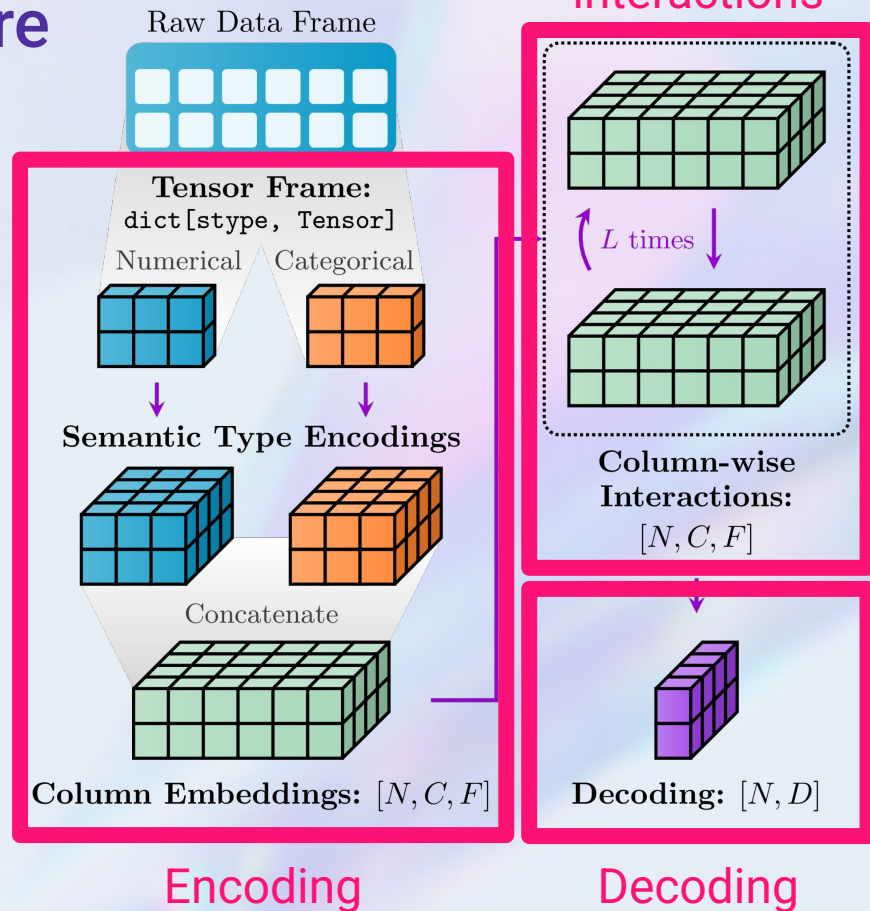
- weighted sum of column embeddings
- MLP over the flattened column embeddings
- ...



Three-stage Model Architecture

Many models fit within three-stage framework

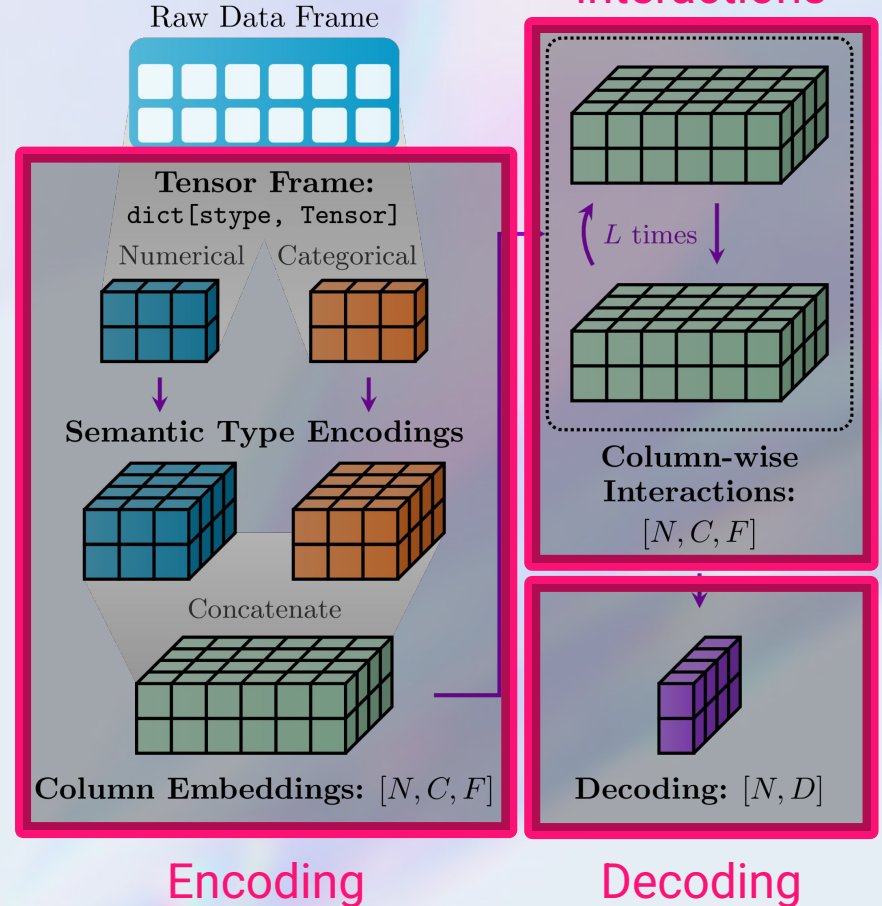
| Model | Trompt | ExcelFormer | FT-Transformer |
|-------------------------|-------------------|-----------------|-------------------------|
| Encoder | Any | CatBoostEncoder | Any |
| Column-wise interaction | Trompt Cell | Transformer | CLS token + Transformer |
| Decoder | Trompt Downstream | MLP | MLP |



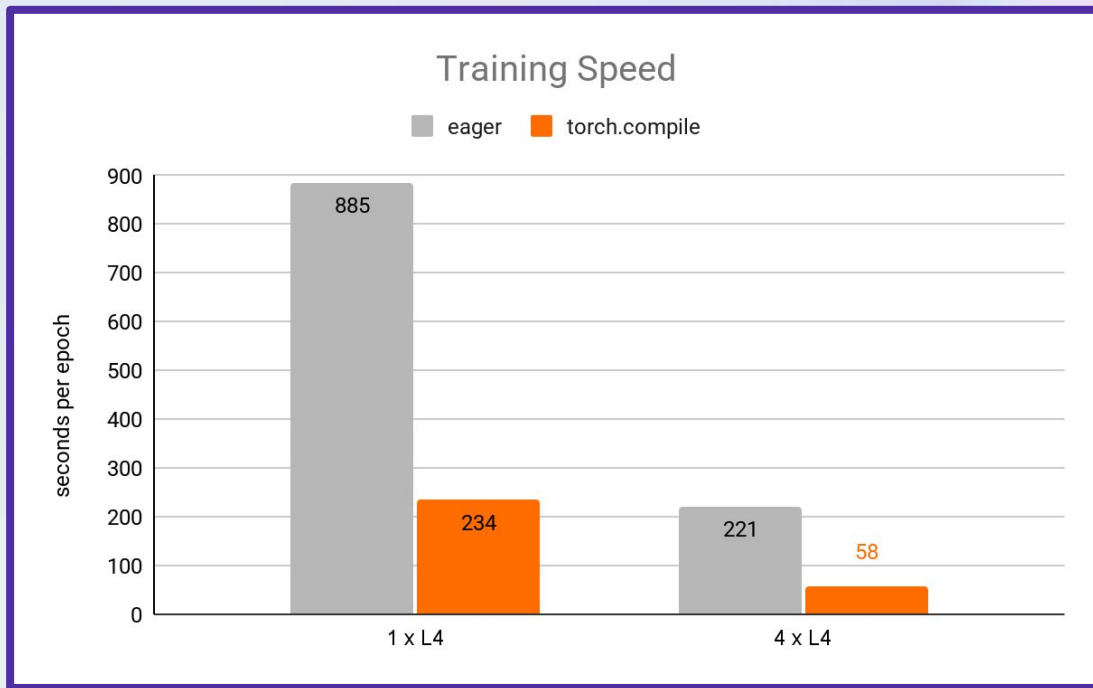
```

1 class MyModel(torch.nn.Module):
2     def __init__(self):
3         self.encoder = StypeWiseFeatureEncoder(
4             out_channels=channels,
5             stype_encoder_dict={
6                 stype.categorical: EmbeddingEncoder(),
7                 stype.numerical: LinearEncoder(),
8             },
9         )
10        self.convs = torch.nn.ModuleList([
11            TabTransformerConv(
12                channels=channels,
13                num_heads=num_heads,
14            ) for _ in range(num_layers)
15        ])
16        self.decoder = torch.nn.Linear(
17            channels,
18            out_channels,
19        )
20
21    def forward(self, tf: TensorFrame) -> Tensor:
22        x, _ = self.encoder(tf)
23        for conv in self.convs:
24            x = conv(x)
25        return self.decoder(x.mean(dim=1))

```



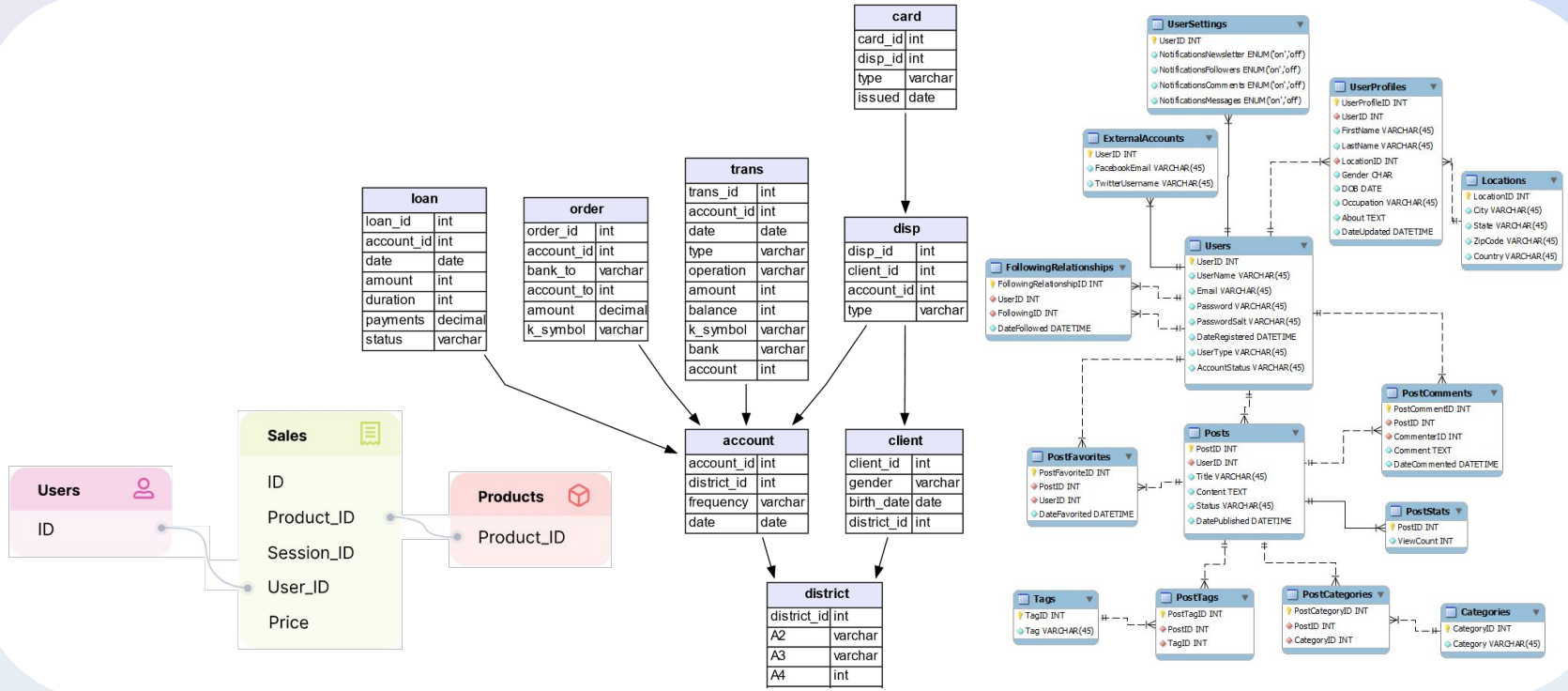
Compatible with torch.compile & DDP!



```
$ python examples/trompt.py --dataset Higgs --batch_size 512  
$ python examples/trompt.py --dataset Higgs --batch_size 512 --compile  
$ python examples/trompt_multi_gpu.py --dataset Higgs --batch_size 512  
$ python examples/trompt_multi_gpu.py --dataset Higgs --batch_size 512 --compile
```

Single Table to Multiple Tables – Relational Deep Learning

Relational Data in Real World



E-Commerce

Finance

Social Media



Relational Data Is a Graph

A **node** denotes a row in a table

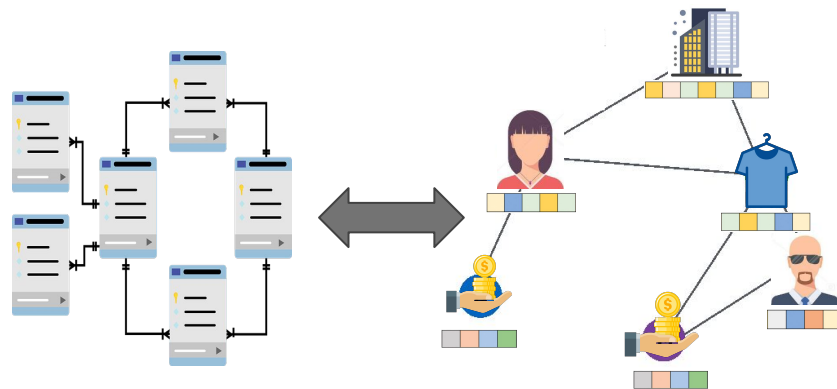
- A user in `USERS` table
- A product in `Sales` table

An **edge** denotes a pkey-fkey relationship between nodes

- `Users.ID` ↔ `Sales.User_ID`
- `Products.Product_ID` ↔ `Sales.Product_ID`

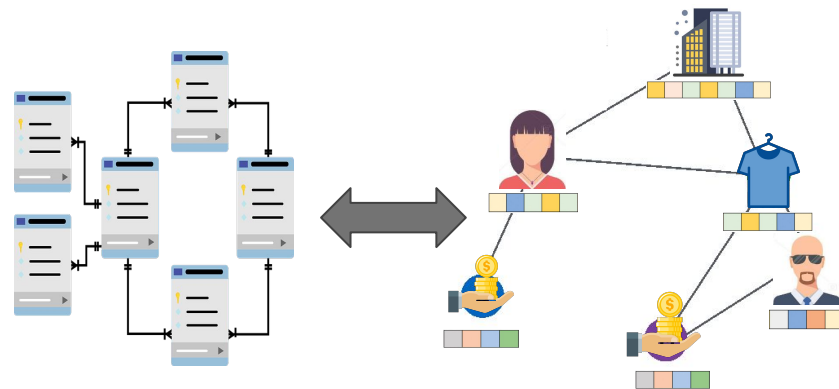
Nodes features are column features:

- User features: age, post code, signup date
- Product features: price, category tags, image, rating



Graph ML Tasks on Relational Data

| Problem | Type |
|---|------------------------------------|
| <u>Age Prediction</u> | Regression |
| <u>Prediction of Movie Ratings</u> | Regression |
| <u>Lifetime Value</u> | Temporal Regression |
| <u>Active Purchasing Customer LTV</u> | Temporal Regression |
| <u>Active Purchasing Customer Churn</u> | Temporal Binary Classification |
| <u>Fraud Detection</u> | Temporal Binary Classification |
| <u>Next Item Category Prediction</u> | Temporal Multiclass Classification |
| <u>Probability of Liking an Item</u> | Non-Temporal Classification |
| <u>Item to Item Similarity</u> | Static Link Prediction |
| <u>Top 25 Most Likely Purchases</u> | Temporal Link Prediction |
| <u>Item Recommendation</u> | Temporal Link Prediction |
| <u>Top 25 "High Value" Purchases</u> | Temporal Classification/Ranking |

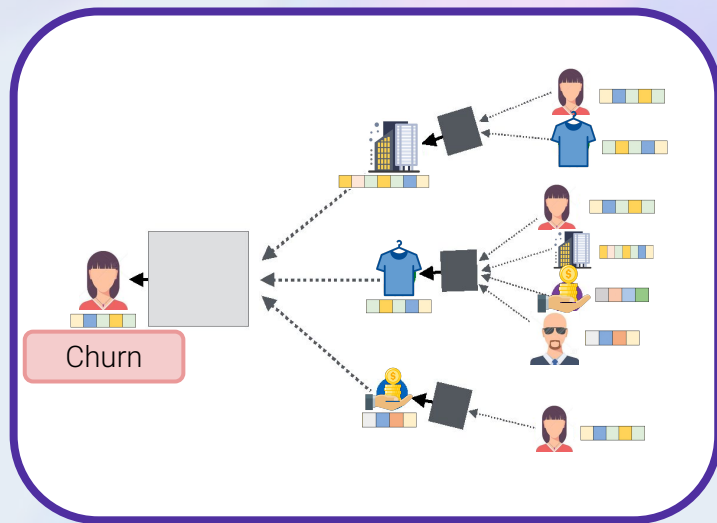


<https://kumo.ai/docs/examples/predictive-query>




Graph Neural Networks on Relational Data

GNNs aggregate information from tables to make predictions

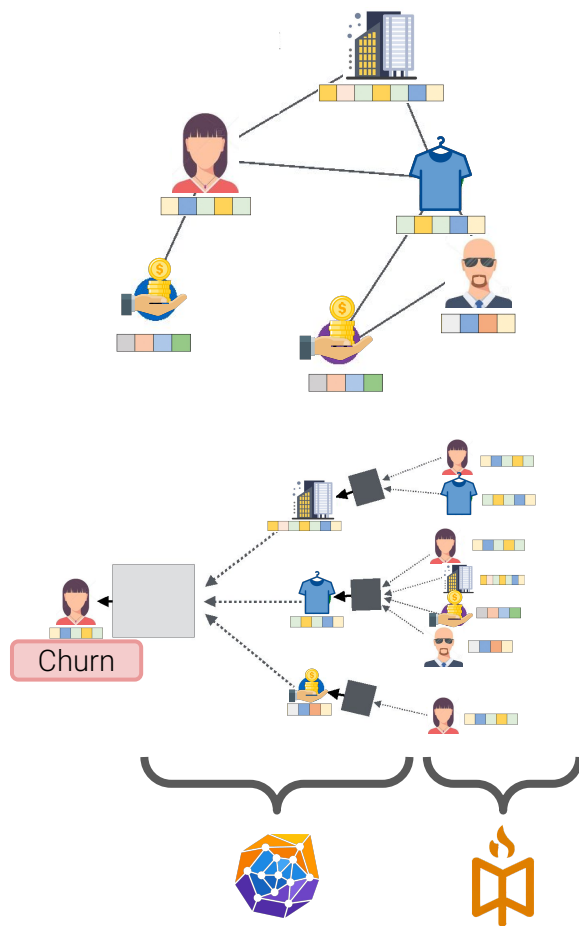
- ✓ No feature engineering
- ✓ No temporal information leakage
via temporal sampling
- ✓ More accurate model
via lossless graph representation



Tabular Models and GNNs

-  Step 1: Sample a subgraph
-  Step 2: Compute node embeddings for nodes in the subgraph
-  Step 3: Perform message passing across table with GNNs

Optimize tabular model and GNNs end-to-end for your task



```

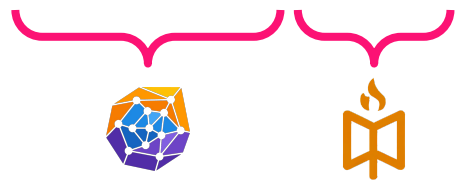
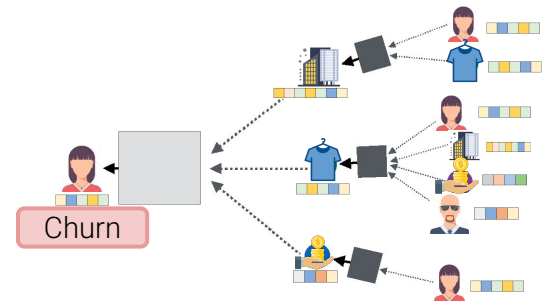
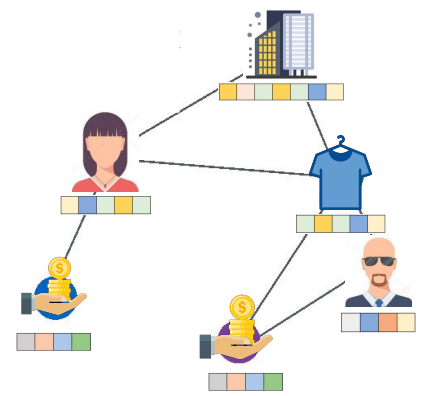
1 class MyGNN(torch.nn.Module):
2     def __init__(self):
3         self.encoder = torch.nn.ModuleDict({
4             "Users": MyTabularModel(...),
5             "Products": Trompt(...),
6         })
7         self.gnn = to_hetero(GraphSAGE(...))
8
9     def forward(
10        self,
11        tf_dict: dict[str, TensorFrame],
12        edge_index_dict: dict[str, Tensor],
13    ) -> Tensor:
14        x_dict = {
15            self.encoder[table_name][tf]
16            for table_name, tf in tf_dict.items()
17        }
18        return self.gnn(x_dict, edge_index_dict)

```




RelBench: A Benchmark for Deep Learning on Relational Databases
<https://arxiv.org/abs/2407.20060>
<https://github.com/snap-stanford/re1bench>

ContextGNN: Beyond Two-Tower Recommendation Systems
<https://arxiv.org/abs/2411.19513>
<https://github.com/kumo-ai/ContextGNN>



PyTorch Frame Community & Summary

- ✓ Support *for* multi-modal features
- ✓ Support *for* LLMs and foundation models
- ✓ Support *for* Relational Deep Learning

 /pyg-team/pytorch-frame

 /pyg-team/pytorch_geometric

```
pip install torch-geometric
pip install pytorch-frame
```

What's coming?

- Add more layers
- Add more models
- Integrate NestedTensor?
- Integrate fbgemm kernels?

Next steps:

- Run examples
- Use your own data
- Build your own models
- Check out RelBench
- Check out ContextGNN
- Contribute!
- Join PyG Slack!

https://join.slack.com/t/torchgeometric/shared_invite/zt-2ppirqs10-_krkJVMqCeXiP92jDCRbig

Community Contributors

anas-rz
berkekisin
Borda
crunai
DamianSzwichtenberg
drivanov
eliazonta
February24-Lee
HoustonJ2013
itsjayway
jyansir
kaidic
NeelKondapalli
puririshi98
rishabh-ranjan
SimonPop
toenshoff
xnuohz

Core Contributors

akihironitta
rusty1s
vid-koci
weihua916
wsad1
XinweiHe
yiweny
zechengz