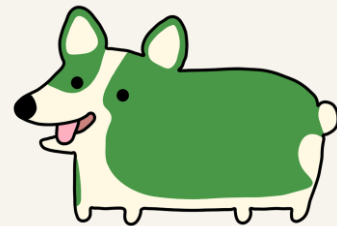
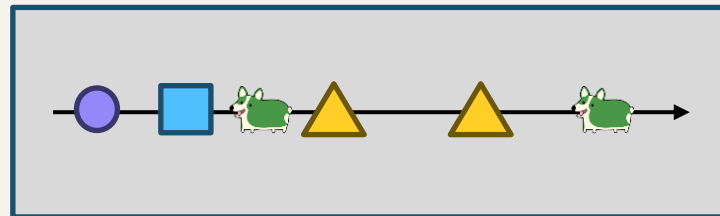


# Universal representations for event sequences: financial transactional data and beyond



Alexey Zaytsev

Assistant professor  
LARSS laboratory,  
Skoltech

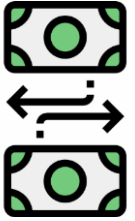
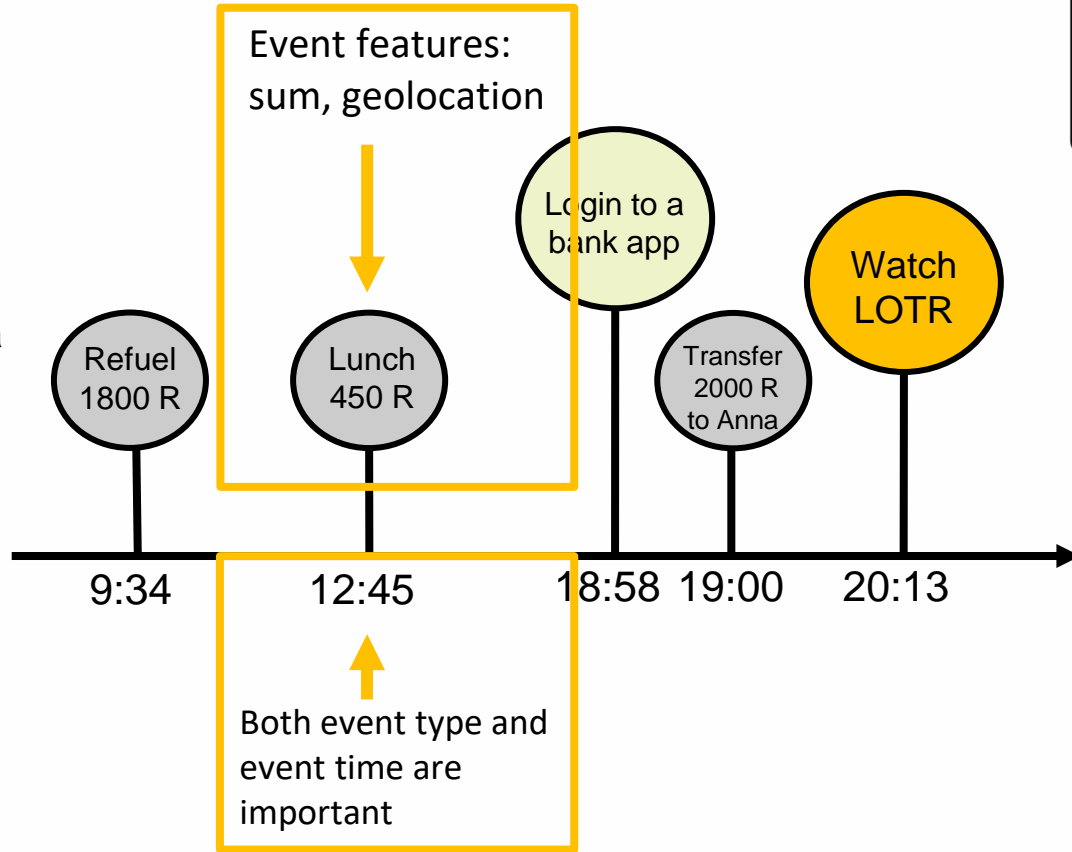


# Event sequences: intro

# Event sequences data



Alex, 35, man  
works in academia



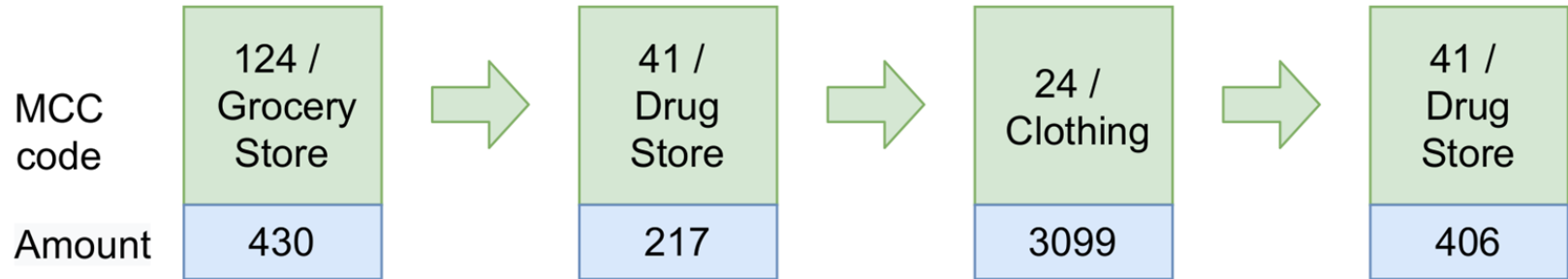
# Discrete sequential transactional data

Transaction records data sequence includes:

- MCC (Merchant Category Codes)
- Purchase amount
- Time values
- Transaction location
- ...

Data characteristics:

- Heterogeneous features
- Non-regularity of observations
- Varying lengths of sequences



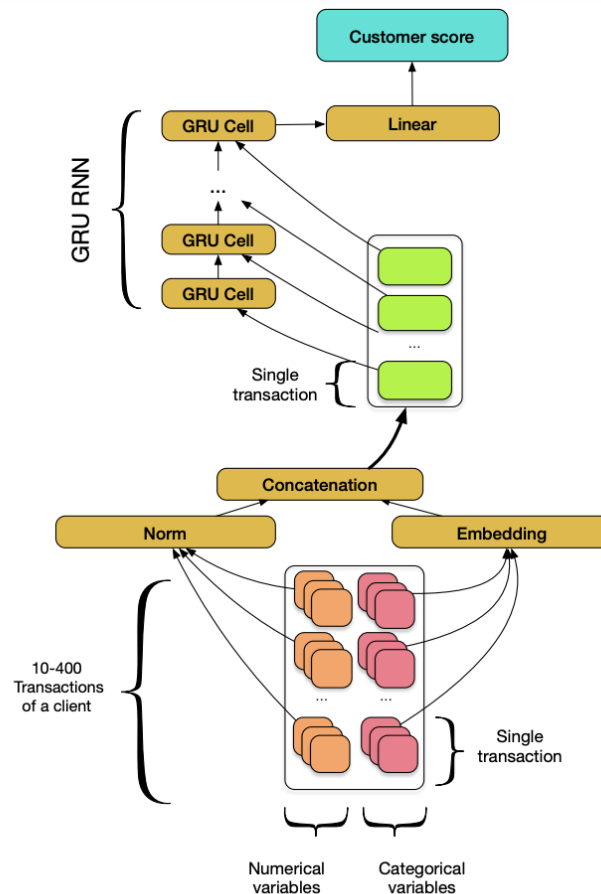
# One way: Supervised approach

Library pytorch-lifestream

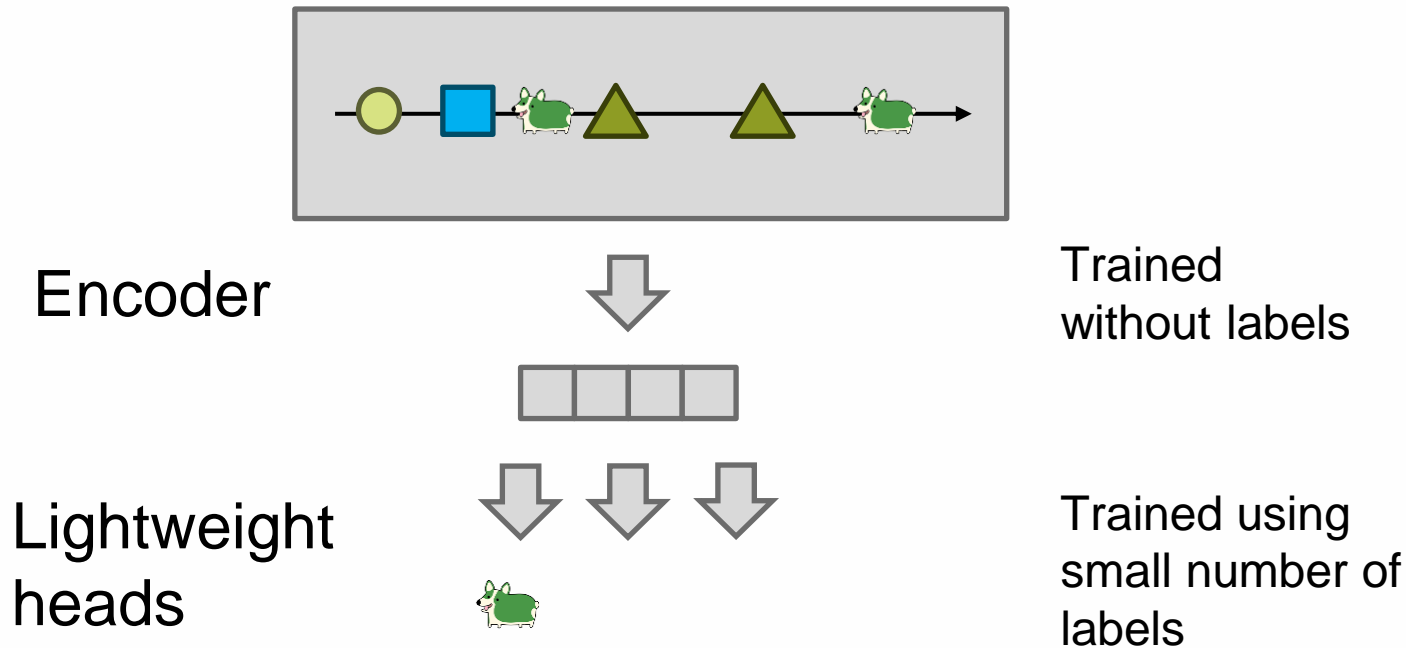
Recurrent (or Transformer) Neural Network with self-supervised contrastive learning

	ROC AUC	N Features
<b>Logistic regression</b>	0.78	~ 400
<b>LGBM</b>	0.81	~ 7000
<b>E.T.-RNN</b>	0.83	12

- Requires labeled data!

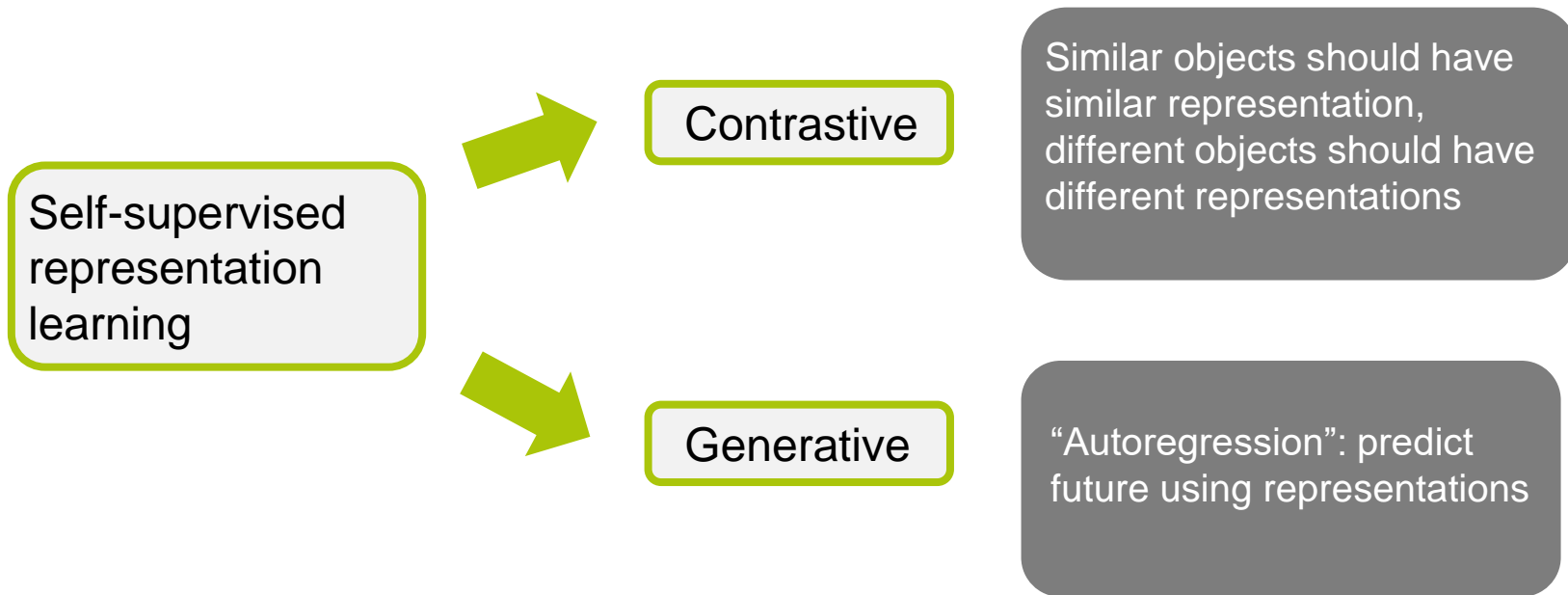


# Our way: self-supervised approach



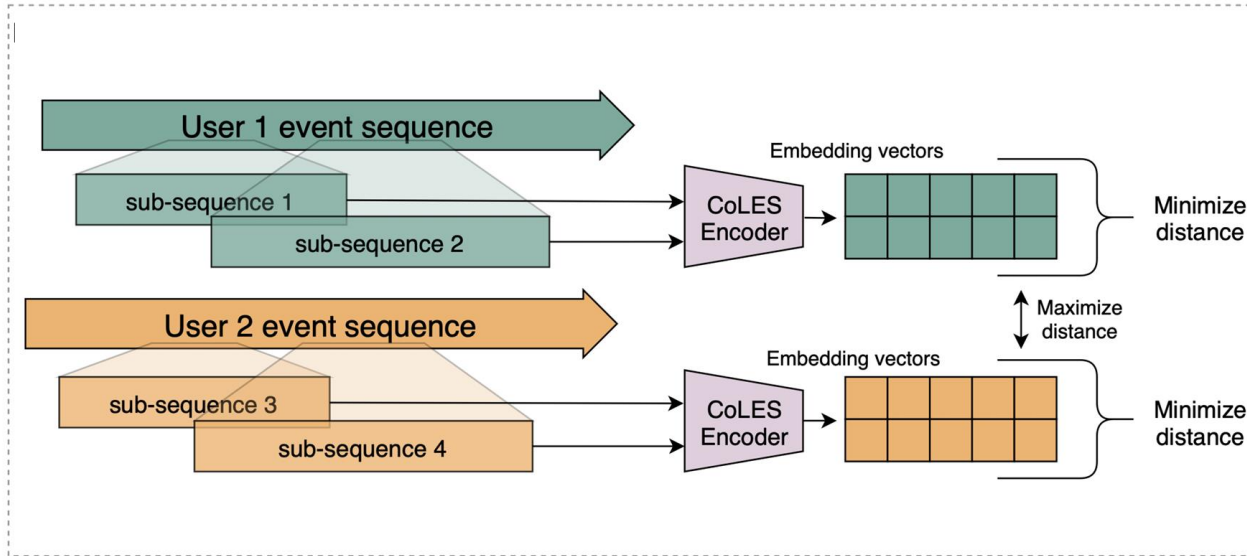
# Types of self-supervised learning

# Two main ideas of self-supervised learning: generative and contrastive





# CoLES contrastive learning



# Contrastive learning for sequential data

Weak augmentation: jitter-and-scale strategy

Strong augmentation: permutation-and-jitter strategy

$$\mathcal{L} = \lambda_1 \cdot (\mathcal{L}_{TC}^s + \mathcal{L}_{TC}^w) + \lambda_2 \cdot \mathcal{L}_{CC}$$

Predict future representation  
from the current context

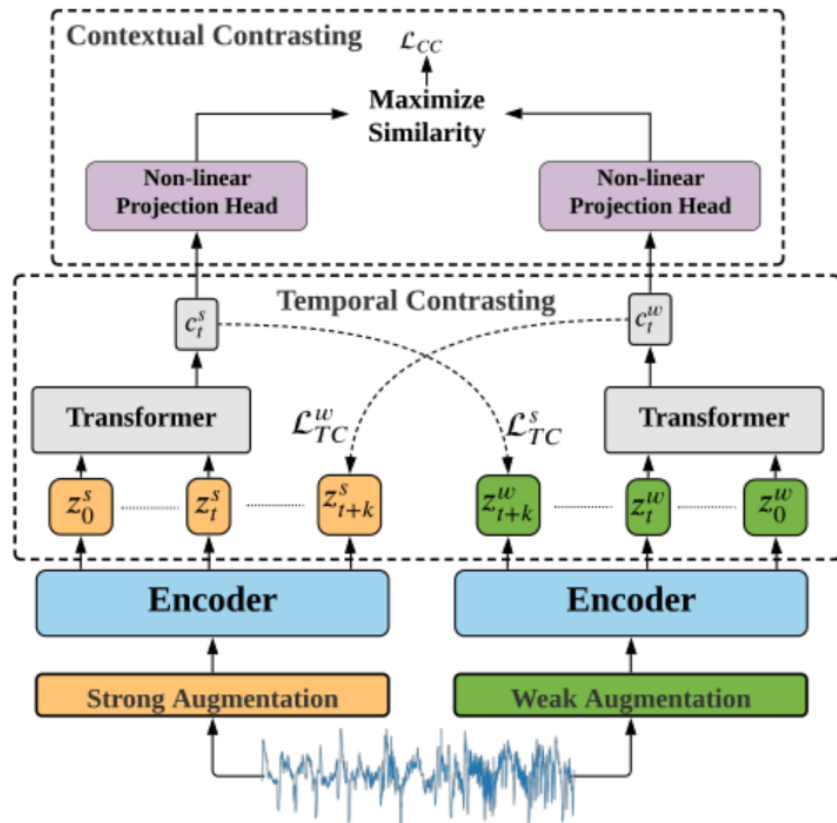
Compare  
contexts

Context:

$$c_t = f_{AR}(Z_{\leq t}),$$

Representation:

$$z_t = f_{enc}(x_t)$$

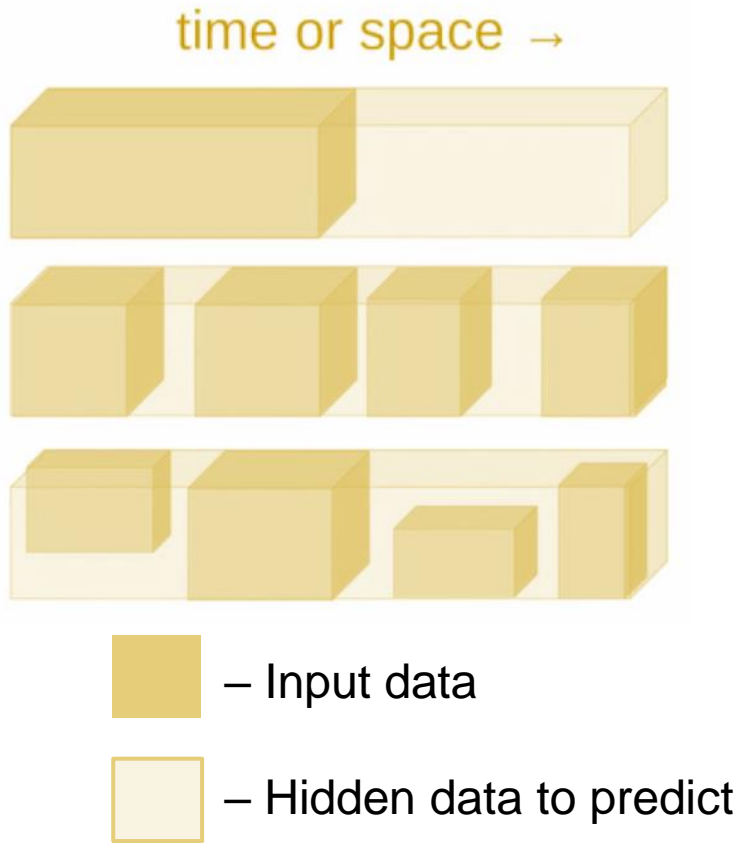


# Generative models: masking

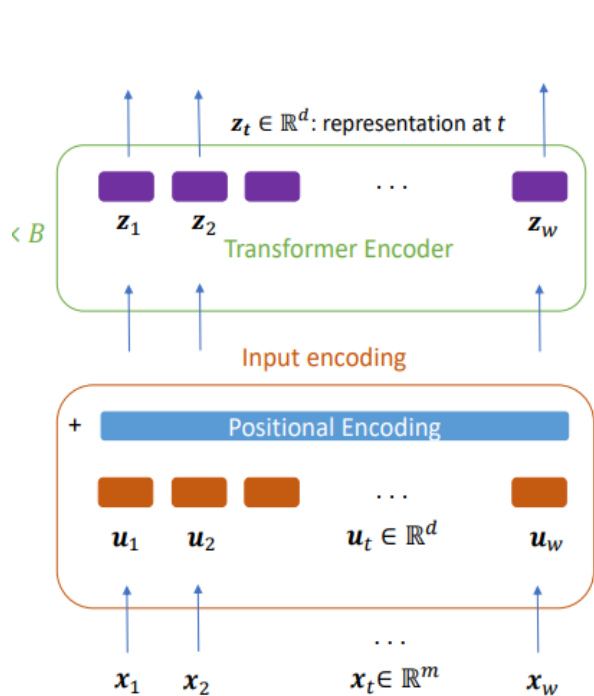
Training steps:

1. Hide some part of the data
2. Try to recover it via representation learning

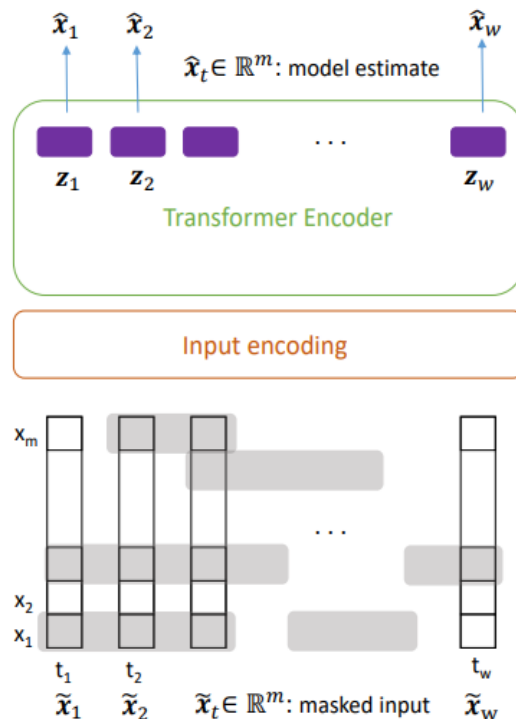
- A. Predict the future from the past
- B. Predict the invisible from the visible
- C. Predict occluded, masked or corrupted part



# Time-series unsupervised representations



Time series encoding via Transformers



Masking for model training

Zerveas, George, et al. A transformer-based framework for multivariate time series representation learning. KDD. 2021.

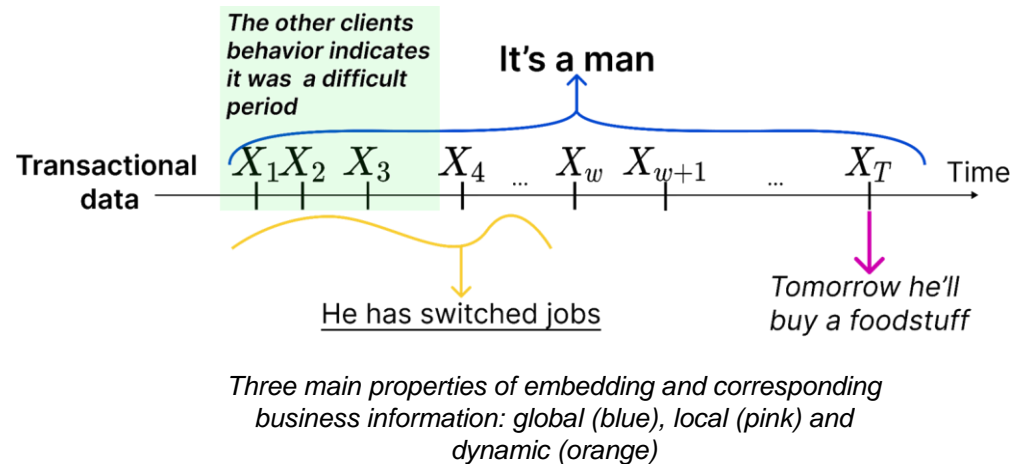
# **Desired properties of embeddings**

# Properties of event sequence embeddings

**Goal: to obtain a good encoder for transactional data**

Three main properties of local embedding for transactional data:

1. **Global property** - describe a client in general;
2. **Local property** – describe a client's state at a particular moment in time;
3. **Dynamic property** - the embeddings should change with time, reflecting the changes in the client's behavior.



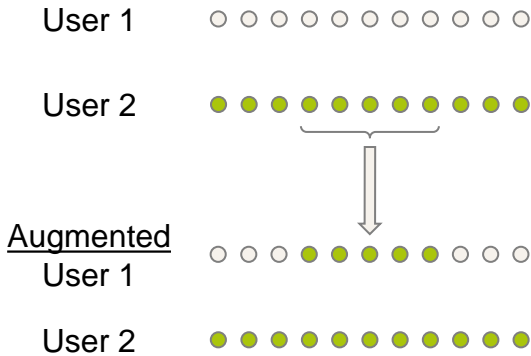
# Global and local quality of the models

- *Global validation* – solve a downstream task via a boosting model, get ROC AUC;
- *Local validation* – two approaches:
  - a. predict the next event type (MCC) via MLP, get ROC AUC instead of likelihood;
  - b. predict a local downstream target (churn/default state at the moment) via MLP, estimate ROC AUC.

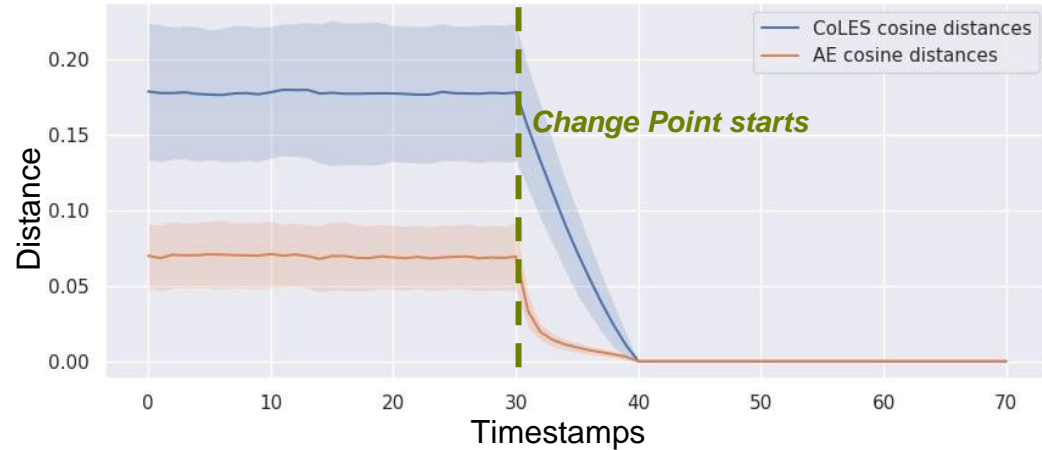
# CoLES (contrastive) vs AE (generative): reaction to change

We also evaluate the models' **ability to detect user behavior change**. See an artificial change.

**Experiment:** "A poor man won a lottery".



*Augmentation procedure. User 1 transactions were replaced with User 2 transactions. We compare User 2 to the augmented User 1.*



*Cosine distance between embeddings obtained from raw users and augmented ones. Snapshot near the Change Point*

We expect embedding during the “augmented” area will be close to each other and far during other timestamps.



# Global properties of models

Ranks for a local problem

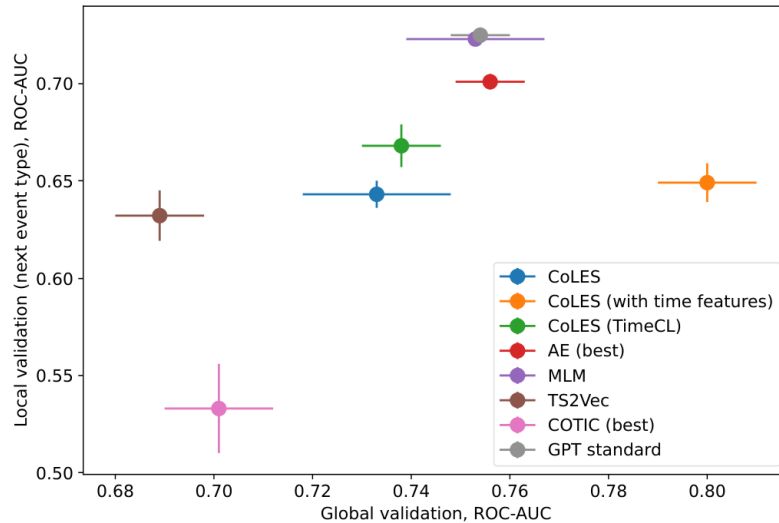
	Age	Churn	Default	HSBC	Mean
AR <sup>§</sup>	1	1	1	1	1.00
MLM <sup>§</sup>	3	1	2	1	1.75
CoLES ext. <sup>†</sup>	3	2	2	3	2.50
AE <sup>§</sup>	2	3	4	2	2.75
CoLES <sup>†</sup>	3	4	3	3	3.25
Best baseline	4	4	5	3	4.00
TS2Vec <sup>†</sup>	6	4	5	4	4.75
A-NHP <sup>‡</sup>	5	5	6	4	5.00
NHP <sup>‡</sup>	5	5	6	4	5.00
COTIC <sup>‡</sup>	6	6	7	5	6.00

Ranks for a global problem

	Age	Churn	Default	HSBC	Mean
CoLES ext. <sup>†</sup>	1	1	1	1	1.00
CoLES <sup>†</sup>	1	3	1	1	1.50
MLM <sup>§</sup>	2	2	2	2	2.00
Best baseline	1	4	2	2	2.25
AR <sup>§</sup>	3	2	1	3	2.25
AE <sup>§</sup>	4	2	2	2	2.50
NHP <sup>‡</sup>	5	2	2	3	3.00
COTIC <sup>‡</sup>	6	3	1	4	3.50
TS2Vec <sup>†</sup>	2	5	2	5	3.50
A-NHP <sup>‡</sup>	5	3	3	4	3.75

Models are colour-coded: **blue** for generative, **green** for contrastive and **fuchsia** for TPP.

# Comparison of local and global properties of models



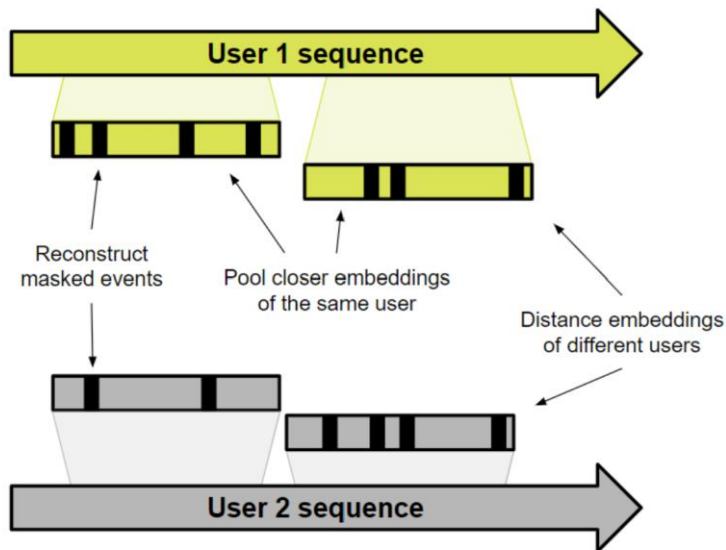
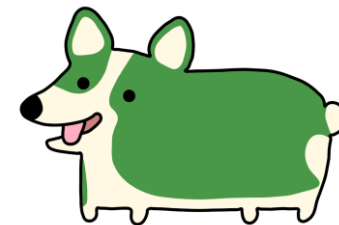
## Main conclusions:

- GPT is better in local task.
- CoLES with time features is a clear leader in global validation.

Bazarova, Alexandra, et al. "Universal representations for financial transactional data: embracing local, global, and external contexts." *arXiv preprint arXiv:2404.02047* (2024).

# **Combining contrastive learning and autoregression**

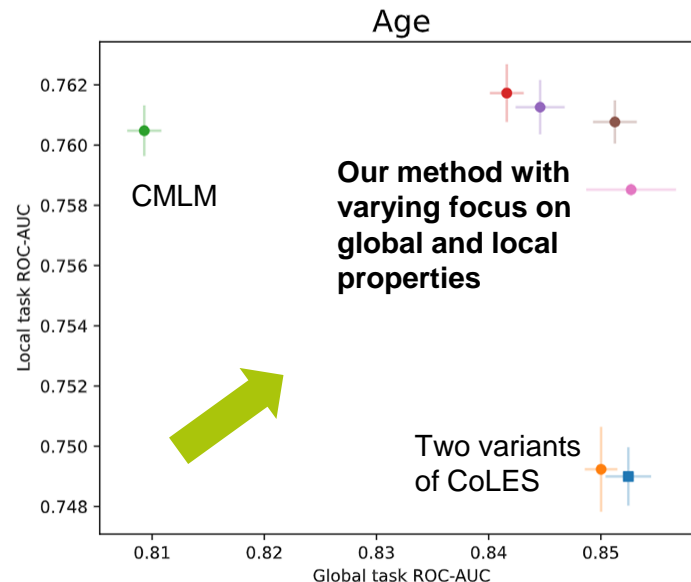
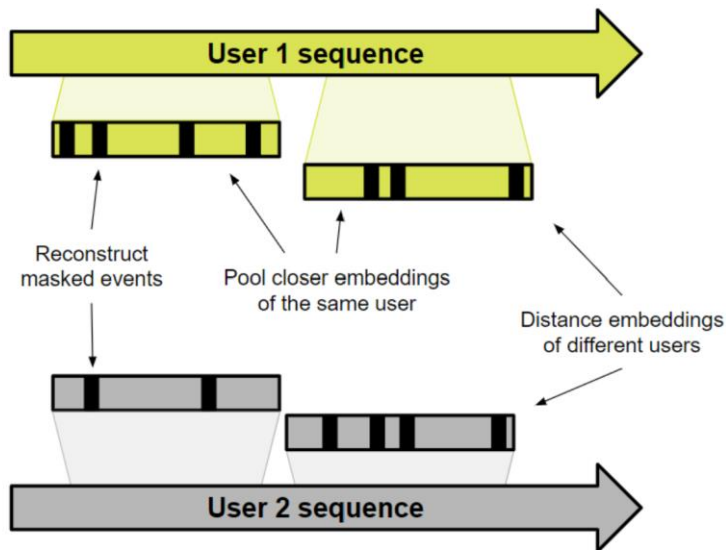
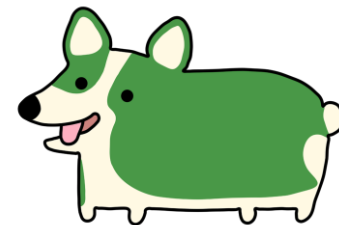
# Combining local and global properties



1. Generative reconstruction embeddings of masked events
2. Contrastive comparison of embeddings from different users

We simultaneously reconstruct embeddings with our CMLM and contrast in CoLES style

# Combining local and global properties

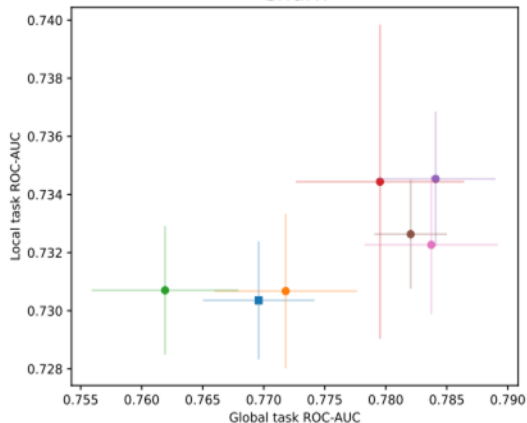


We simultaneously reconstruct embeddings with our CMLM and contrast in CoLES style

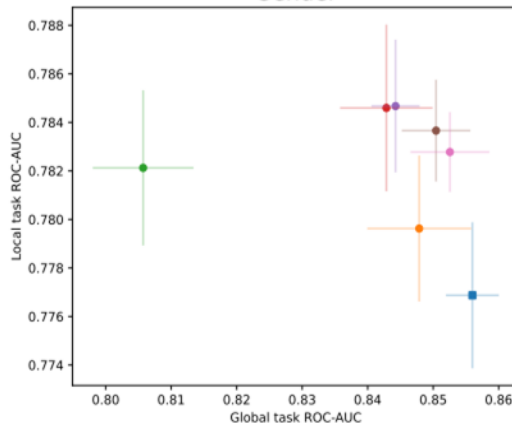
# Results



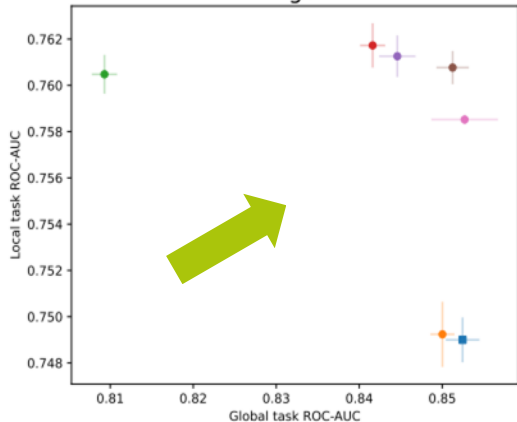
Churn



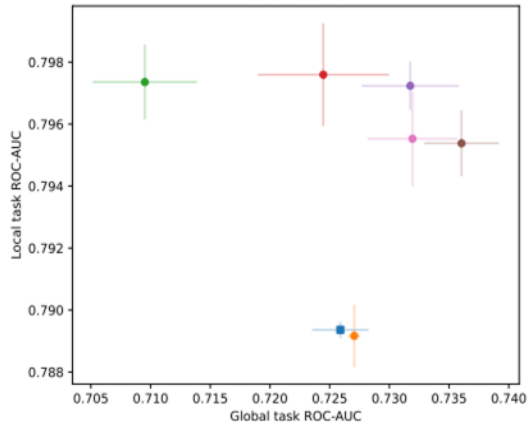
Gender



Age



DataFusion



Method	Global Task ROC-AUC	Local Task ROC-AUC
--------	---------------------	--------------------

**Churn**

CoLES	0.770±0.007	0.730±0.003
CoLES (masking)	0.772±0.007	0.731±0.003
CMLM	0.762±0.010	0.731±0.004
CMLM+CoLES ( $\lambda = 0.1$ )	0.780±0.008	0.734±0.006
CMLM+CoLES ( $\lambda = 0.05$ )	<b>0.784±0.008</b>	<b>0.735±0.004</b>
CMLM+CoLES ( $\lambda = 0.01$ )	0.782±0.005	0.733±0.003
CMLM+CoLES ( $\lambda = 0.005$ )	<b>0.784±0.009</b>	0.732±0.004

**Gender**

CoLES	<b>0.856±0.005</b>	0.777±0.004
CoLES (masking)	0.848±0.009	0.780±0.003
CMLM	0.806±0.009	0.782±0.004
CMLM+CoLES ( $\lambda = 0.1$ )	0.843±0.007	<b>0.785±0.004</b>
CMLM+CoLES ( $\lambda = 0.05$ )	0.844±0.004	<b>0.785±0.003</b>
CMLM+CoLES ( $\lambda = 0.01$ )	0.850±0.005	0.784±0.002
CMLM+CoLES ( $\lambda = 0.005$ )	<u>0.853±0.008</u>	0.783±0.002

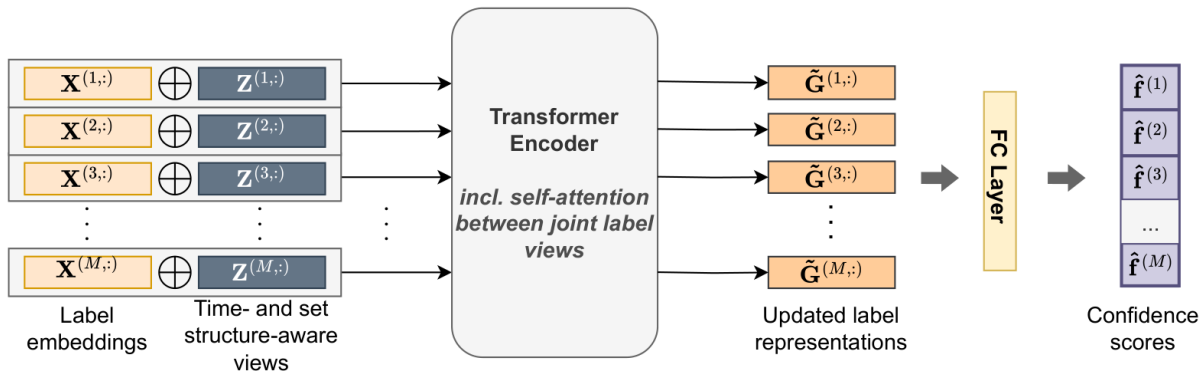
**Age**

CoLES	<u>0.852±0.002</u>	0.749±0.001
CoLES (masking)	0.850±0.001	0.749±0.001
CMLM	0.809±0.002	0.760±0.001
CMLM+CoLES ( $\lambda = 0.1$ )	0.842±0.002	<b>0.762±0.001</b>
CMLM+CoLES ( $\lambda = 0.05$ )	0.845±0.002	<u>0.761±0.001</u>
CMLM+CoLES ( $\lambda = 0.01$ )	0.851±0.002	<u>0.761±0.001</u>
CMLM+CoLES ( $\lambda = 0.005$ )	<b>0.853±0.005</b>	0.759±0.000

**DataFusion**

CoLES	0.726±0.003	0.789±0.000
CoLES (masking)	0.727±0.001	0.789±0.001
CMLM	0.710±0.005	0.797±0.001
CMLM+CoLES ( $\lambda = 0.1$ )	0.724±0.005	<b>0.798±0.001</b>
CMLM+CoLES ( $\lambda = 0.05$ )	0.732±0.005	<u>0.797±0.001</u>
CMLM+CoLES ( $\lambda = 0.01$ )	<b>0.736±0.003</b>	0.795±0.001
CMLM+CoLES ( $\lambda = 0.005$ )	<u>0.734±0.005</u>	0.795±0.001

# You were looking at a wrong self-attention?



We compute self-attention over event types and get prediction of next event type, imposing simple aggregation of temporal encodings.

Our LaNET model is now SOTA for the next basked prediction



Kovtun, Elizaveta, et al. Label attention network for sequential multi-label classification: you were looking at a wrong self-attention. ECAI. 2024.

# Few final words



# Conclusion

- Typical SSL approaches focus on different aspects of embedding properties, also demonstrating generative capabilities
- We propose an SSL hybrid approach CMLM+CoLES that achieve notable improvements in both local and global properties of learned representations.
- Generative models for event sequences data are on their way!

Alexandra Bazarova  
Maria Kovaleva  
Ilya Kuleshov  
Evgenia Romanenkova  
Alexander Stepikin  
Alexandr Yugay  
Elizaveta Kovtun  
Galina Boeva  
Andrey Shulga  
Alexey Zaytsev

**Thanks my lab for help with these slides  
and you for your attention!**

**Thanks for your  
attention!**

---

# Backslides

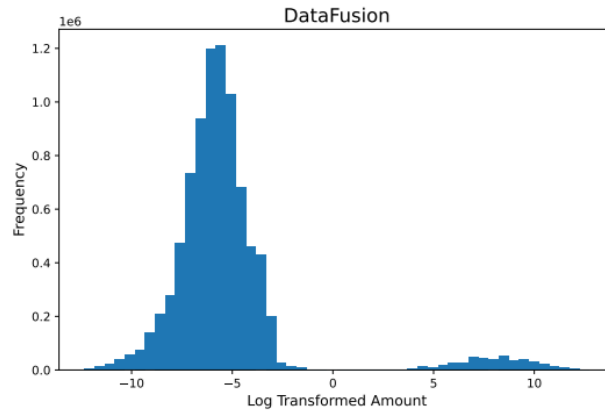
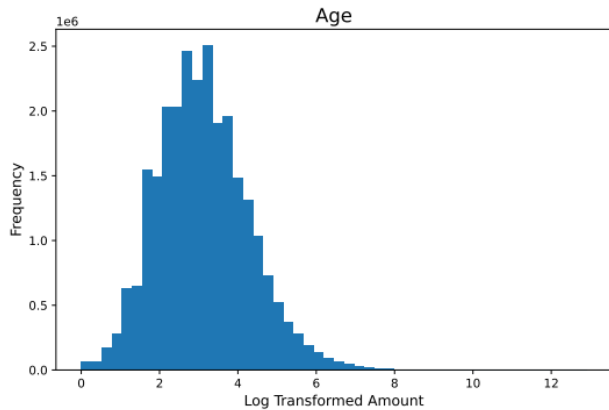
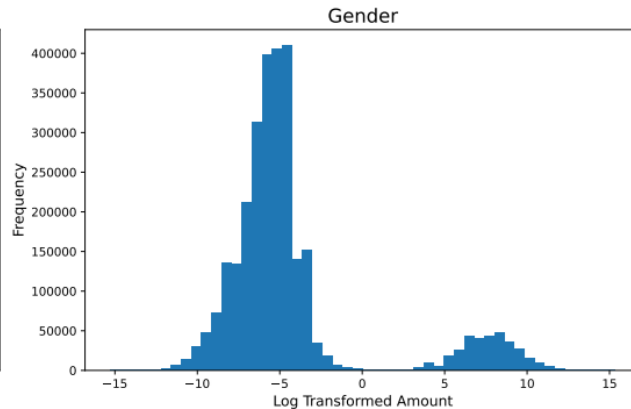
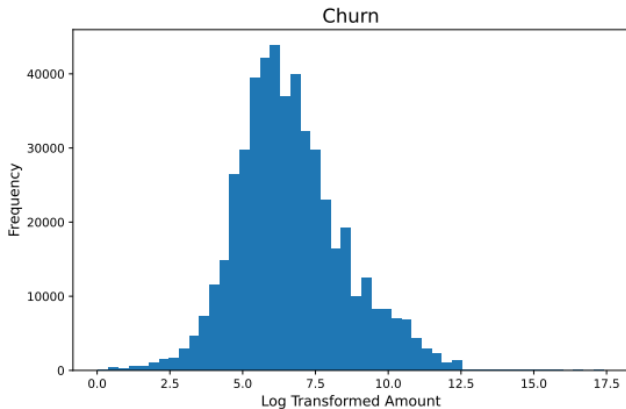


# Experiment design

1. Pretrain models in a self-supervised regime
2. Use the obtained encoder as feature extractor
3. Train another model in a supervised regime on extracted features to solve downstream tasks:
  - Sequence classification
  - Next event type prediction

	<b>Churn</b>	<b>Gender</b>	<b>Age</b>	<b>DataFusion</b>
<b>Num Transactions</b>	490K	2.9M	26M	8.7M
<b>Num Sequences</b>	5K	7.4K	30K	64K
<b>Mean Sequence Length</b>	98.1	388.2	881.7	136.5
<b>Std. Sequence Length</b>	78.1	309.4	124.8	148.9
<b>Num Unique MCC</b>	344	184	202	323

# EDA: Amount



# EDA: Sequence length

