# Agenda

## LUYANG WANG

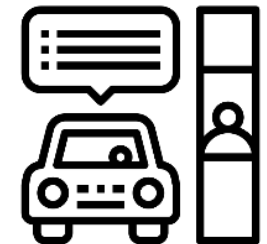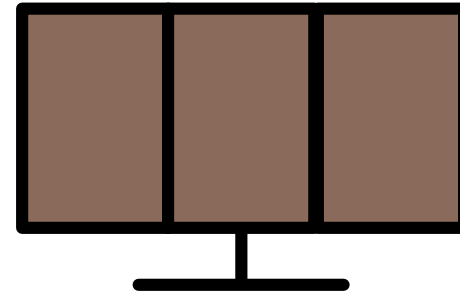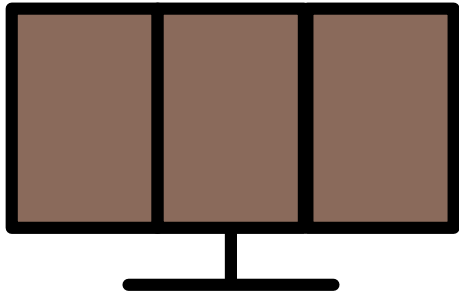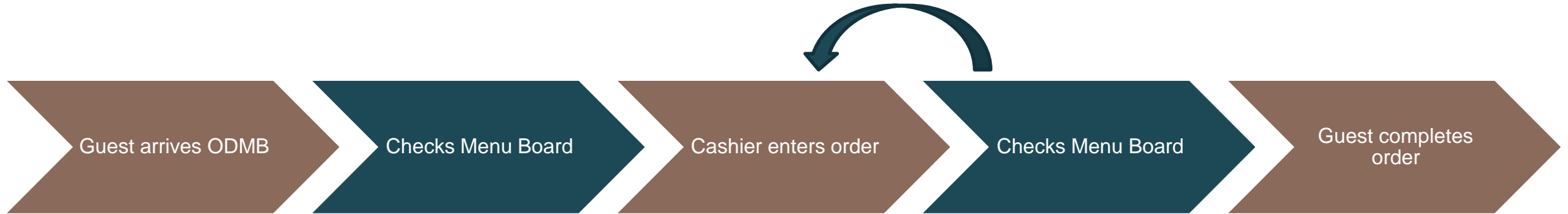- Food recommendation use case

- T$x$T model in detail

## KAI HUANG

- AI on big data

- Distributed training pipeline with Ray on Apache Spark

# Food Recommendation Use Case

# Food Recommendation Use Case

| Guest arrives ODMB | Checks Menu Board | Cashier enters order | Checks Menu Board | Guest completes order |

# Use Case Challenges


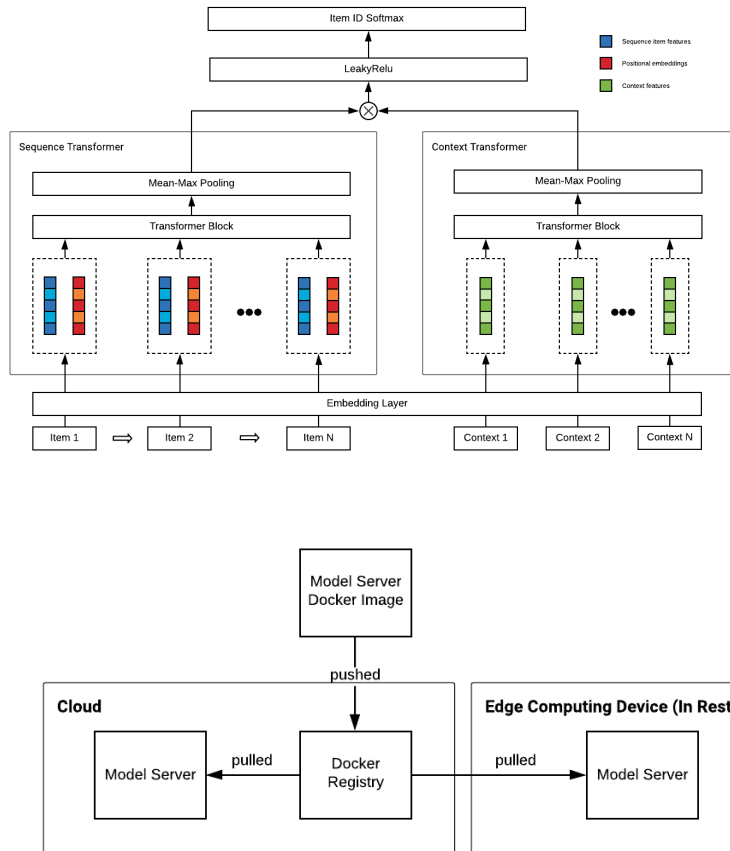
### Challenges

- Lack of user identifiers
- Same session food compatibilities
- Other variables in our use case: locations, weathers, time, etc.
- Deployment challenges

# Use Case Challenges



## Solutions

- Session based recommendation model
- Able to take complex context features into consideration
- Able to be deployed anywhere, both edge / cloud

# Transformer Cross Transformer (T*x*T)

# TxT Model Overview



## Model Components

- ## Sequence Transformer
  - Taking item order sequence as input
- ## Context Transformer
  - Taking multiple context features as input
- ## Latent Cross Joint Training
  - Element-wise product for both transformer outputs

# Model Comparison
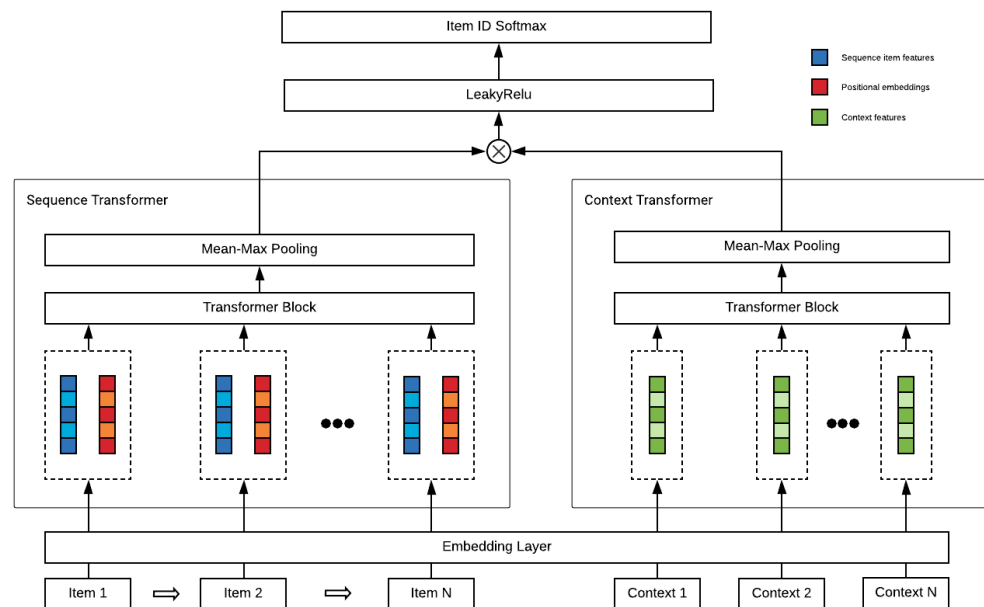
T𝑥T



RNN Latent Cross

# Offline Evaluation

## Offline Training Loss



## Offline Training Result

| Model | Top1 Accuracy | Top3 Accuracy |
|---|---|---|
| RNN | 29.98% | 46.24% |
| Contextual ItemCF | 32.18% | 48.37% |
| RNN Latent Cross | 33.10% | 49.98% |
| TxT | 34.52% | 52.37% |

# Online Performance

## Inference Performance

### Inference Latency (ms)



## A/B Testing Result

| Model | Conversation Rate Gain | Add-on Sales Gain |
|---|---|---|
| RNN Latent Cross (control) | - | - |
| TxT | +7.5% | +4.7% |

# Model Training Architecture

# AI on Big Data

# AI on Big Data

Accelerating Data Analytics + AI Solutions At Scale

- **BigDL: Distributed Deep Learning Framework for Apache Spark**

  https://github.com/intel-analytics/BigDL

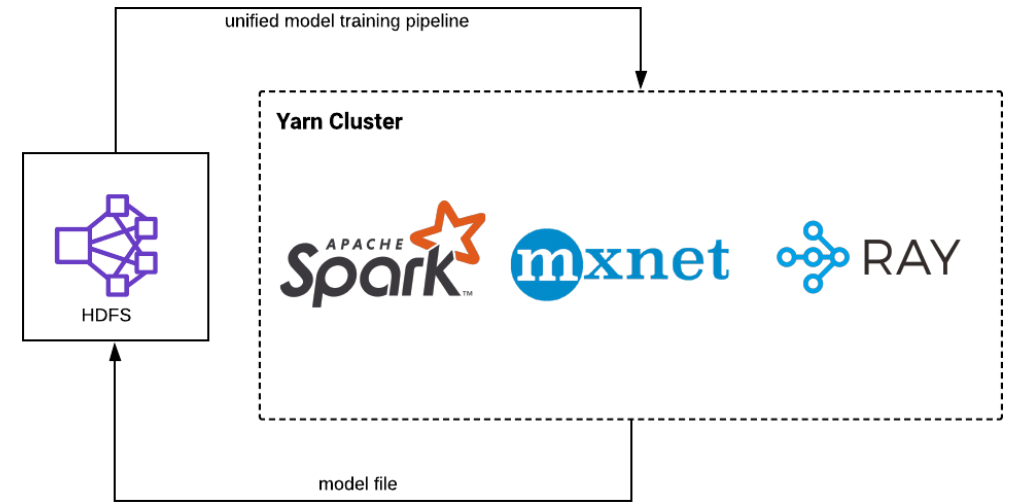- **Analytics Zoo: Distributed TensorFlow, Keras and PyTorch on Apache Spark/Flink & Ray** https://github.com/intel-analytics/analytics-zoo

- We develop *Project Orca* in Analytics Zoo based on Spark and Ray to allow users to easily scale out single node Python notebook across large clusters, by providing:
  - Data-parallel preprocessing for Python AI (supporting common Python libraries such as Pandas, Numpy, PIL, TensorFlow Dataset, PyTorch DataLoader, etc.)
  - Sklearn-style APIs for transparently distributed training and inference (supporting TensorFlow, PyTorch, Keras, MXNet, Horovod, etc.)

  https://github.com/intel-analytics/analytics-zoo/tree/master/pyzoo/zoo/orca

# Ray

Ray is a fast and simple framework for building and running distributed applications.

- Ray Core provides easy Python interface for parallelism by using remote functions and actors.

Ray is packaged with several high-level libraries to accelerate machine learning workloads.

- Tune: Scalable Experiment Execution and Hyperparameter Tuning
- RLlib: Scalable Reinforcement Learning
- RaySGD: Distributed Training Wrappers
- https://github.com/ray-project/ray/

# RayOnSpark

## Seamlessly integrate Ray applications into Spark data processing pipelines.

- Runtime cluster environment preparation.

- Create a *SparkContext* on the drive node and use Spark to perform data cleaning, ETL, and preprocessing tasks.

- *RayContext* on Spark driver launches Ray across the cluster.

- Similar to RaySGD, we implement a lightweight shim layer around native MXNet modules for easy deployment on YARN cluster.

- Each MXNet worker takes the local data partition of Spark RDD or DataFrame from the plasma object store used by Ray.

# End-to-end Distributed Training Pipeline

*Project Orca* provides a user-friendly interface for the pipeline.

- Minimum code changes and learning efforts are needed to scale the training from single node to big data clusters.

- The entire pipeline runs on a single cluster. No extra data transfer needed.

```python
from zoo.orca import init_orca_context
from zoo.orca.learn.mxnet import Estimator

# init_orca_context unifies SparkContext and RayContext
sc = init_orca_context(cluster_mode="yarn", num_nodes, cores, memory)
# Use sc to load data and do data preprocessing.


mxnet_estimator = Estimator(train_config, model=txt, loss=SoftmaxCrossEntropyLoss(),
                            metrics=[mx.metric.Accuracy(), mx.metric.TopKAccuracy(3)])

mxnet_estimator.fit(data=train_rdd, validation_data=val_rdd, epochs=…, batch_size=…)
```

# Conclusion

- Context-Aware Fast Food Recommendation at Burger King with RayOnSpark

https://arxiv.org/abs/2010.06197

https://medium.com/riselab/context-aware-fast-food-recommendation-at-burger-king-with-rayonspark-2e7a6009dd2d

- For more details of RayOnSpark: https://databricks.com/session_na20/running-emerging-ai-applications-on-big-data-platforms-with-ray-on-apache-spark

- More information for Analytics Zoo at:

https://github.com/intel-analytics/analytics-zoo

https://analytics-zoo.github.io/

# Feedback

Your feedback is important to us.

## Don't forget to rate and review the sessions.

**DATA+AI** SUMMIT EUROPE

**#DataTeams #DataAISummit**