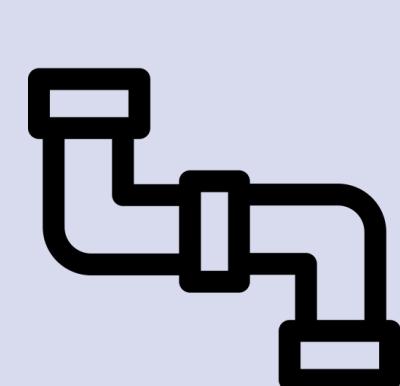


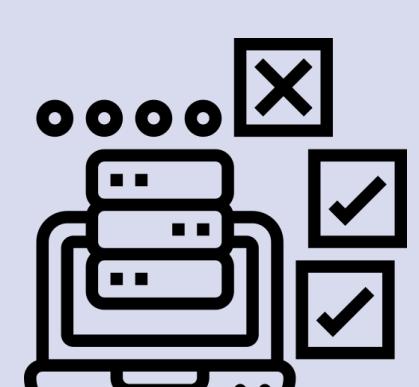
DISTRIBUTED ML Training: A SERVERLESS ARCHITECTURAL APPROACH

Amine Barrak
Fehmi Jaafar (Advisor)
Fabio Petrillo (Co-Advisor)

Serverless Functions & Distributed ML: Challenges, Opportunities, Best Practices



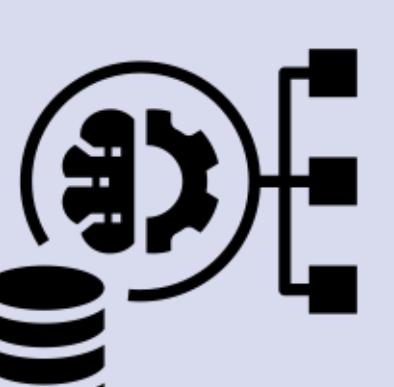
How have serverless functions been utilized in ML pipelines?



How can we secure and ensure fault tolerance in serverless training?



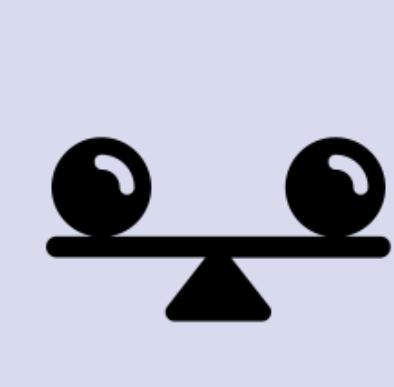
How can serverless functions be used to speed up ML training?



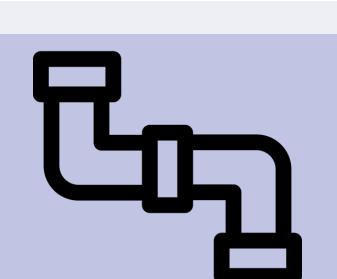
How do we propose a fully serverless ML training architecture?



How can communication overhead be mitigated in serverless distributed training?

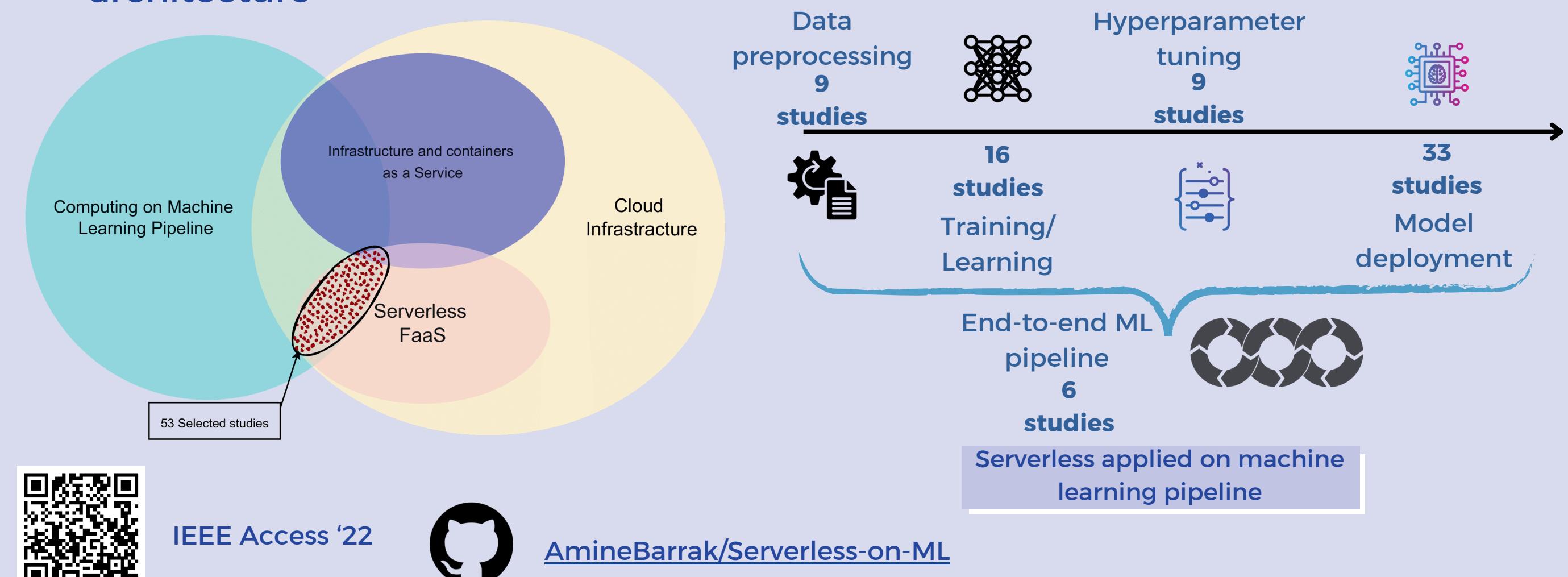


What did we learn from comparing serverless ML frameworks?



Serverless Functions utilization in Machine Learning Pipelines

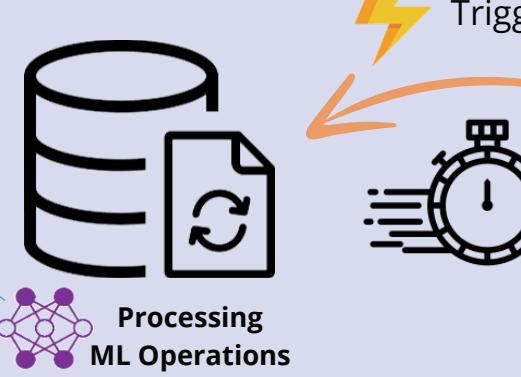
Conducting a systematic mapping study on ML systems applied on serverless architecture



Reduce communication overhead: Perform ML operations within the database

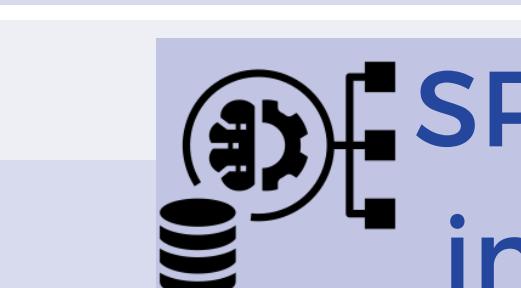
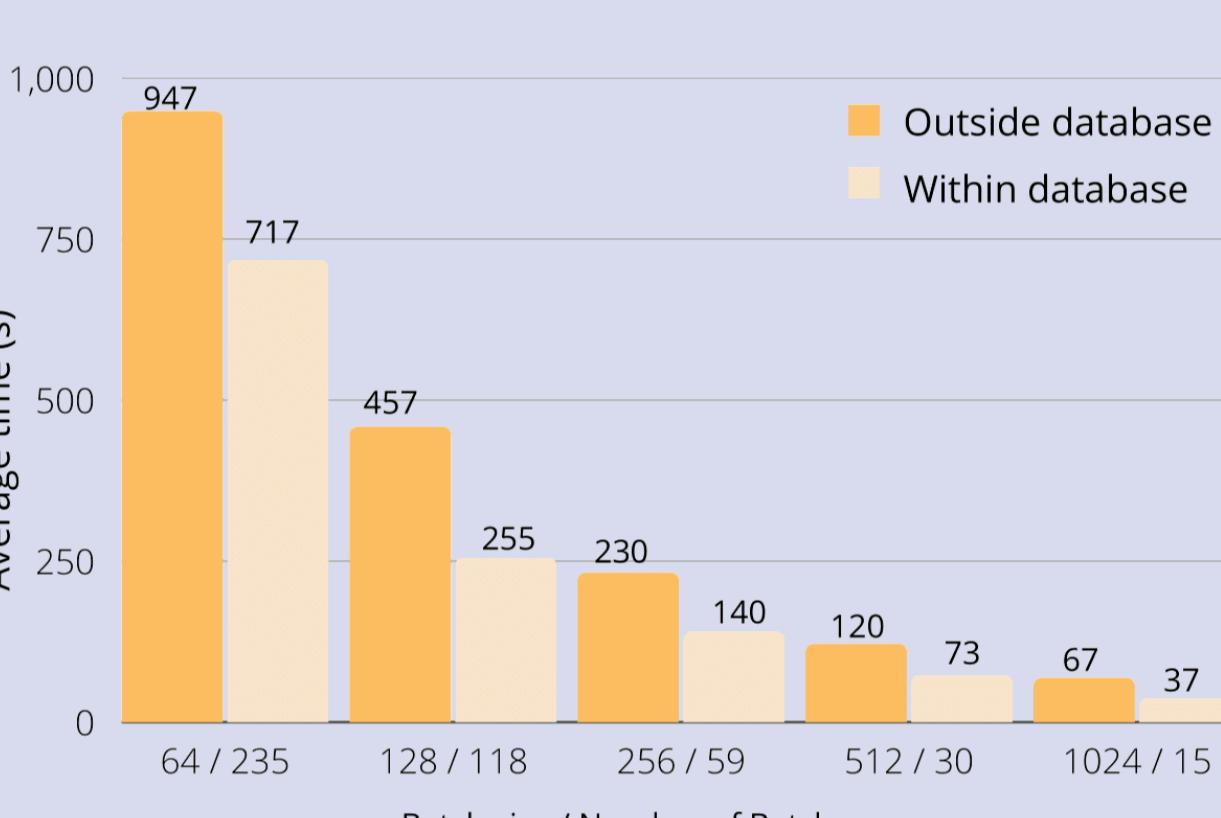


Model update within DB

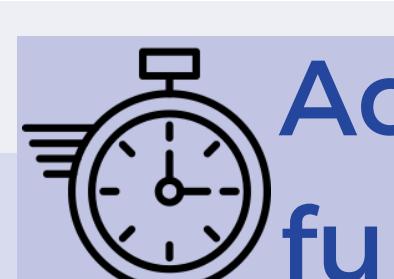
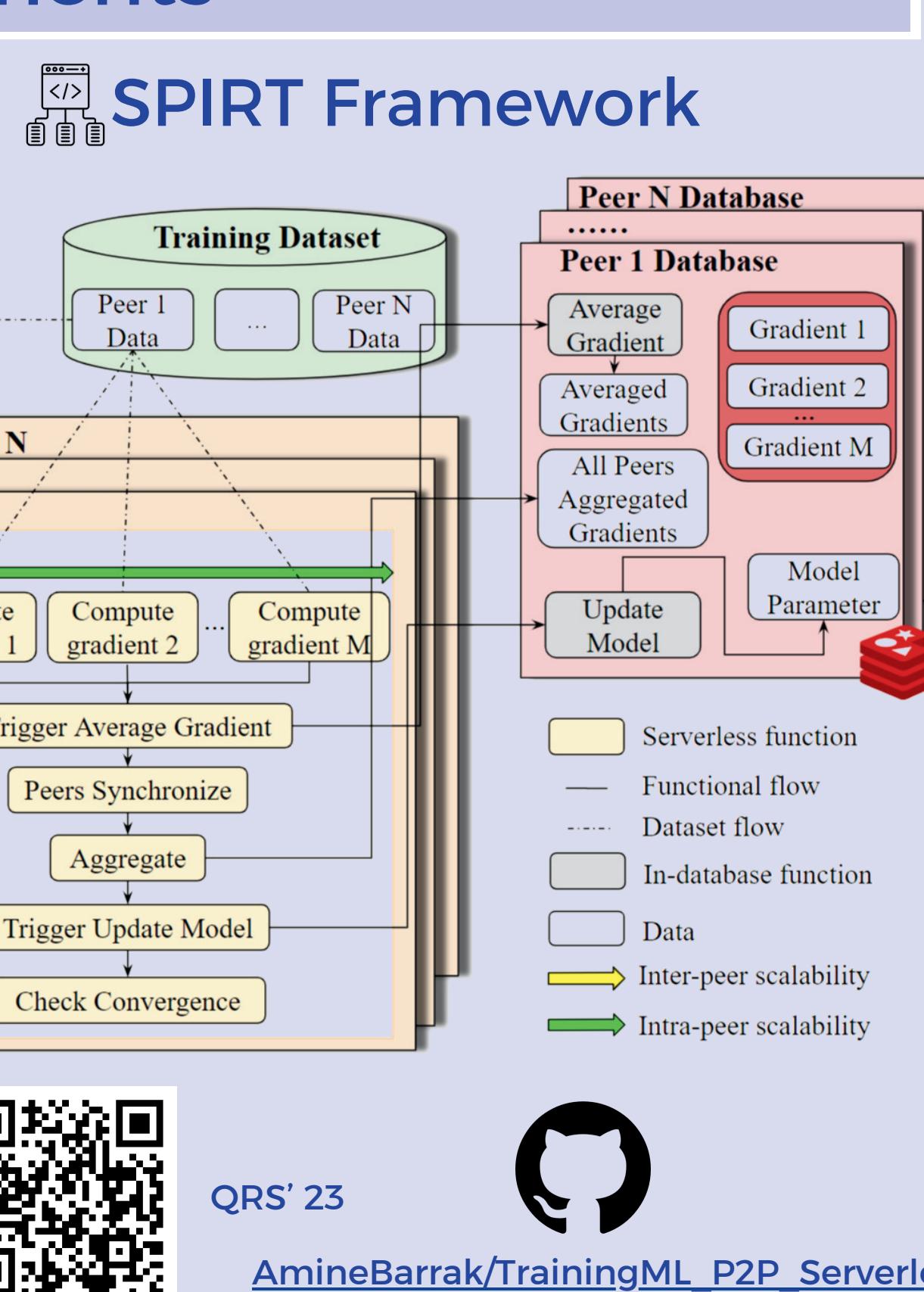
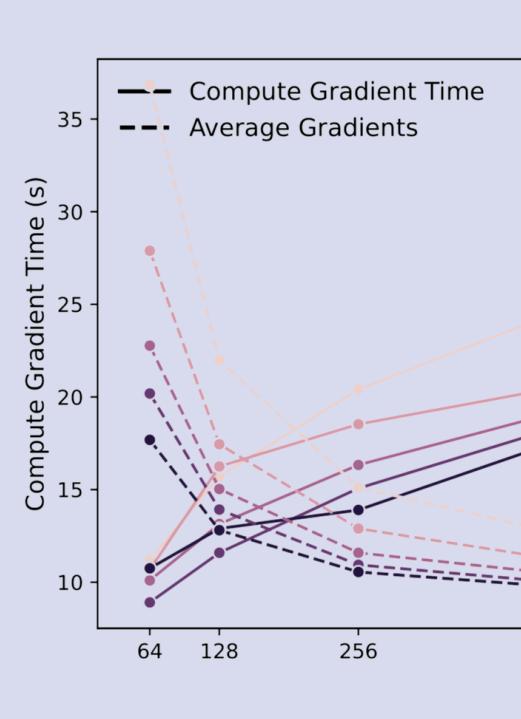
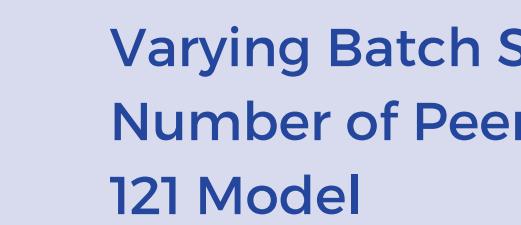


RedisAI's in-database operations achieve a remarkable 82% reduction with ResNet-18

Gradients Averaging within DB



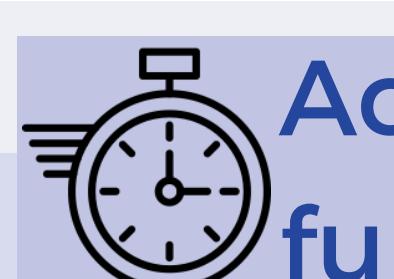
SPIRT : Framework for training ML workflow in serverless environments



Accelerating ML training with serverless functions

Sequential computation

Increased Processing Time



gradients accumulation



Parallel Gradient Computation

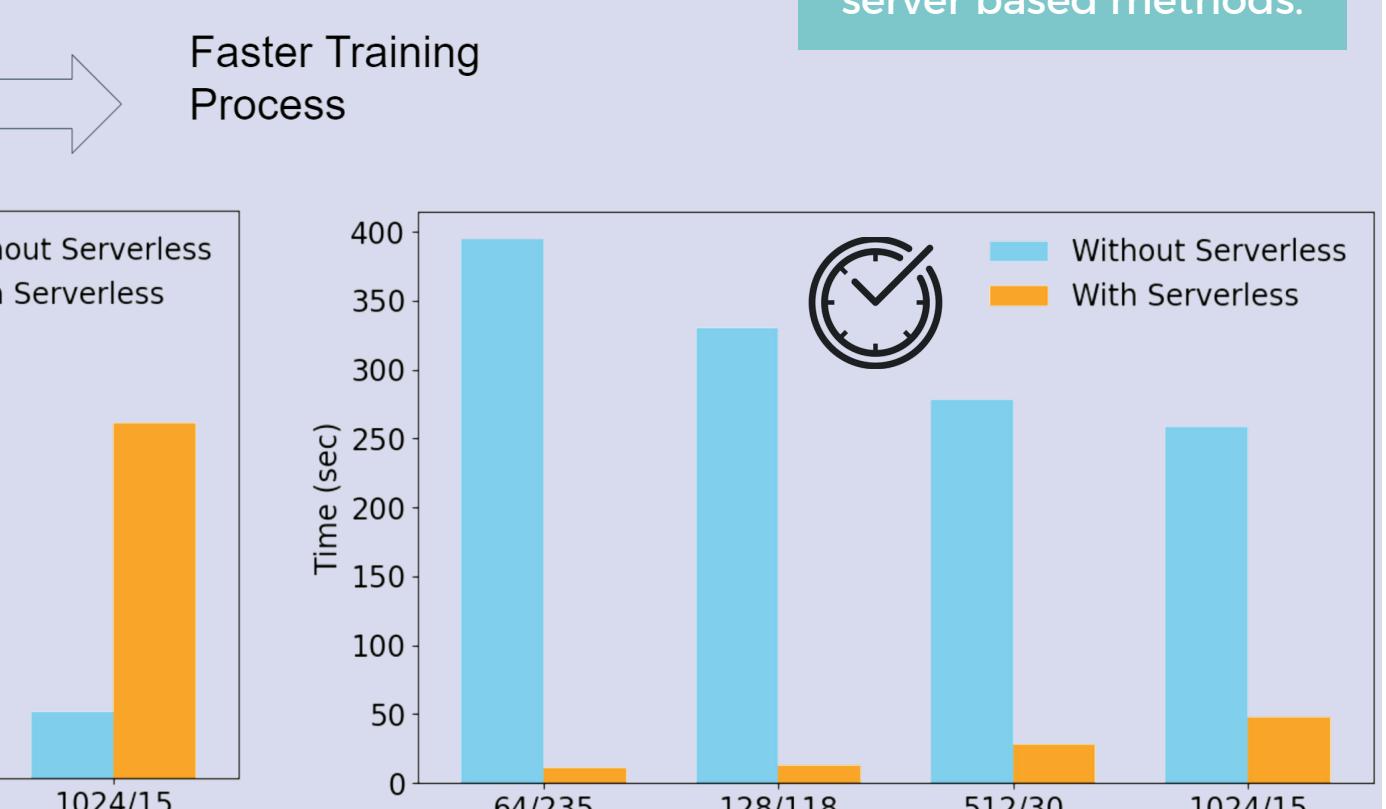


Faster Training Process

Time improvement by 97.34% but cost 5.3 times more than traditional server based methods.

IC2E' 23

AmineBarrak/PeerToPeerServerless

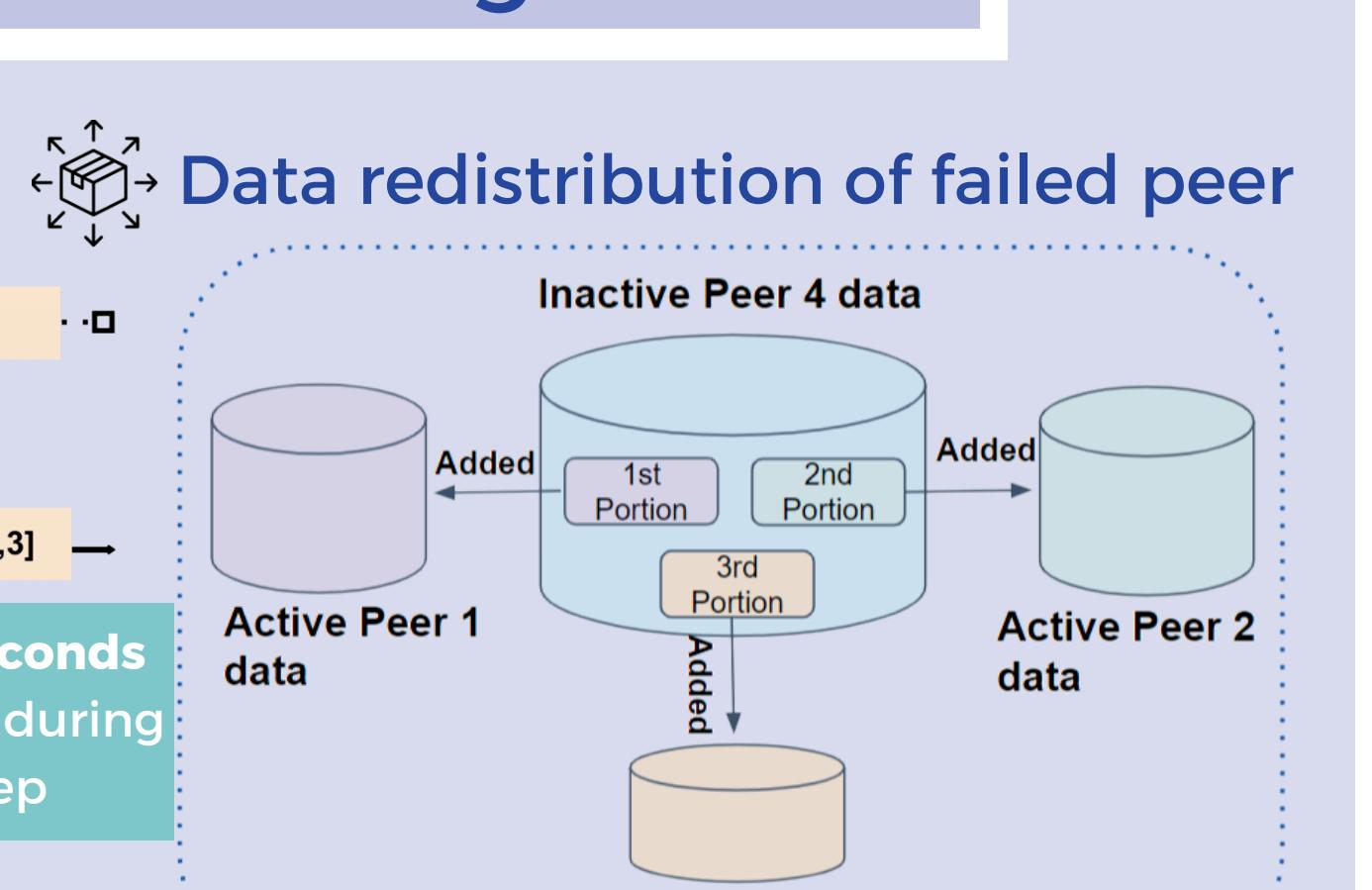


Secure & Fault tolerant communication in serverless distributed training



Adding new peers

Peers authentication mechanism / Asymmetric Encryption



Serverless ML training frameworks comparison

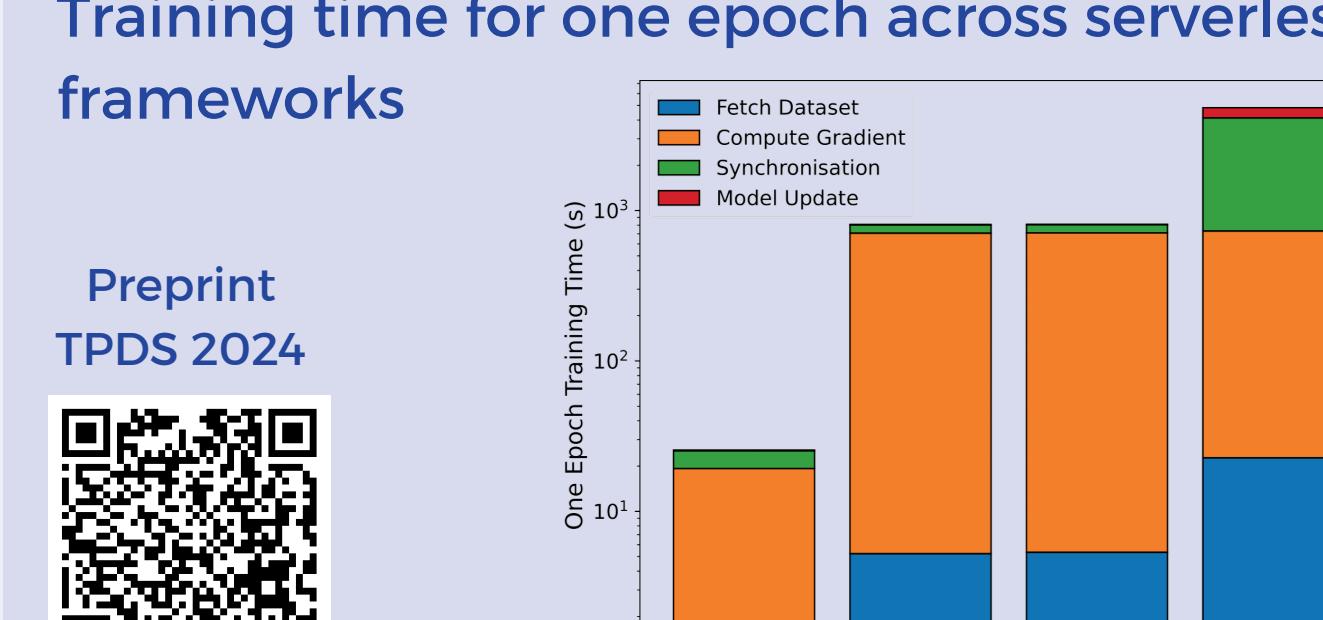
Typical Serverless Training Workflow

- Fetch dataset : Each worker fetches its dataset partition
- Compute gradient : Independent training of workers on local batches/compute gradients
- Synchronization : Gradients upload and aggregation in shared database by workers
- Model Update : Update each local model with the aggregated gradients

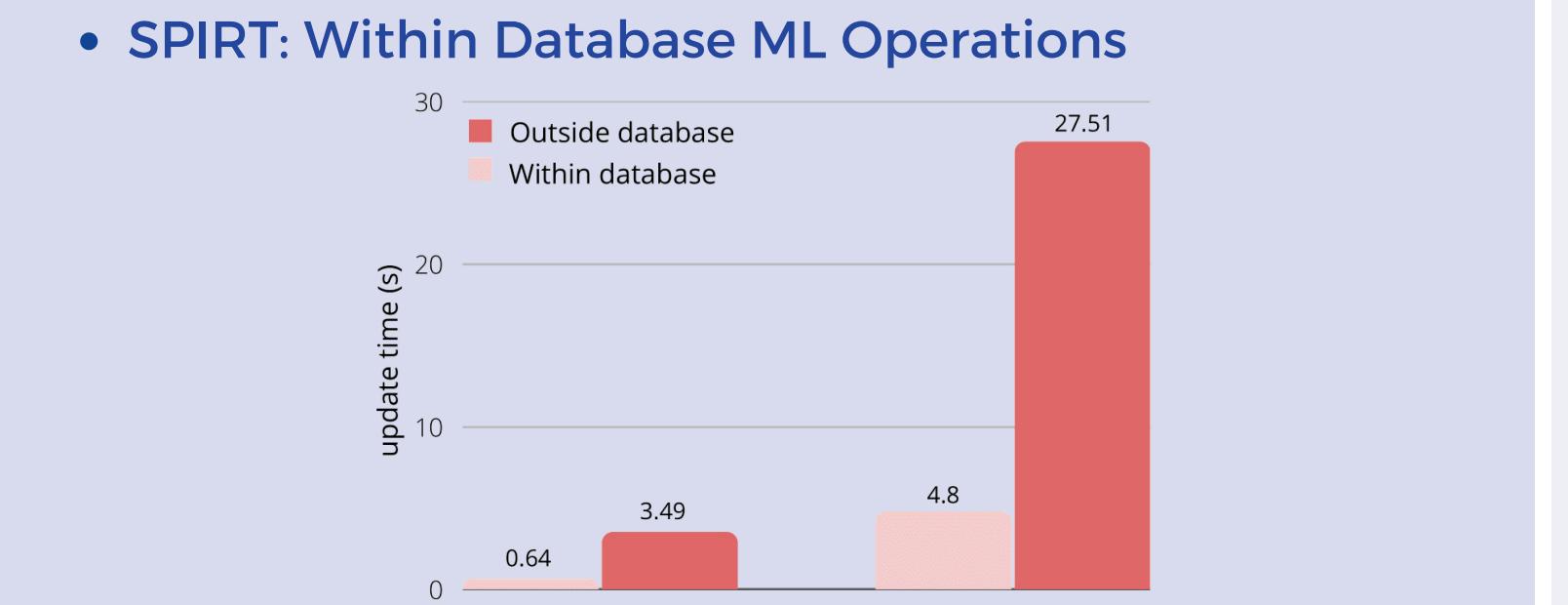
Split the training workflow into separate serverless functions: Optimisation in RAM usage



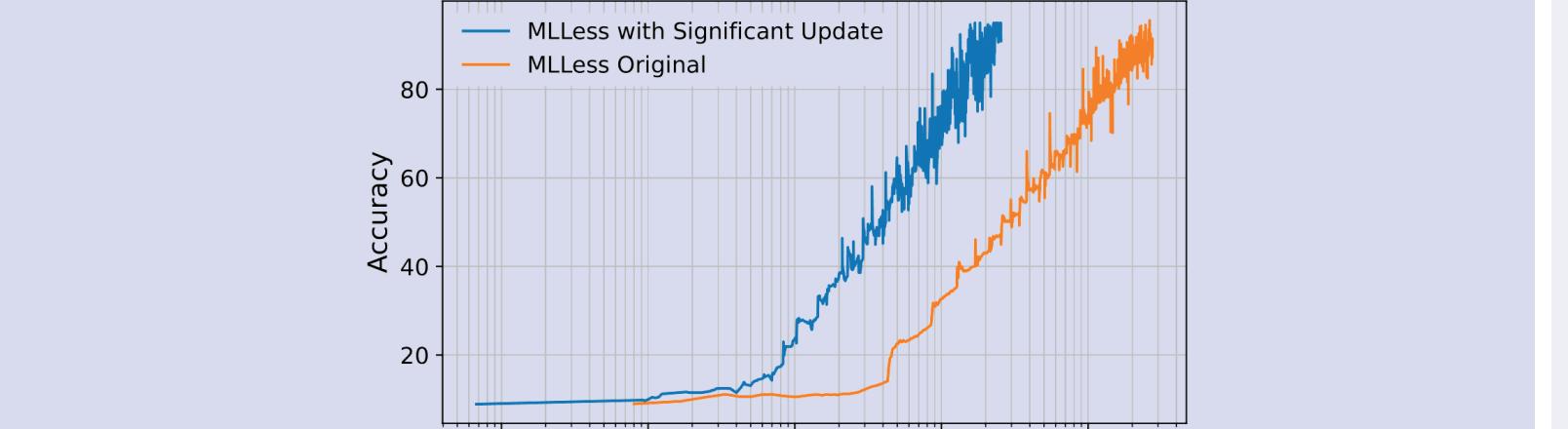
Training time for one epoch across serverless frameworks



Reduce Communication Overhead



MLless: Synchronise only Significant Updates



LambdaML: ScatterReduce Vs. AllReduce

