SNIA | CLOUD STORAGE CSTI | TECHNOLOGIES

Why Distributed Edge Data is the Future of Al

Live Webinar

October 3, 2023

10:00 am PT / 1:00 pm ET

Today's Presenters







Erin Farr Vice Chair SNIA Cloud Storage Technologies Initiative Storage CTO Office, IBM **Rita Wouhaybi** Senior Al Principal Engineer Office of the CTO, Network & Edge Group Intel Heiko Ludwig Senior Manager Al Platforms IBM Research



The SNIA Community

200	2,500	50,000
Corporations,	Active	Worldwide
universities, startups,	contributing	IT end users and
and individuals	members	professionals





What

We

Educate vendors and users on cloud storage, data services and orchestration



Support & promote

business models and architectures: OpenStack, Software Defined Storage, Kubernetes, Object Storage



Understand Hyperscaler requirements Incorporate them into standards and programs



SNIA Legal Notice

- The material contained in this presentation is copyrighted by SNIA unless otherwise noted.
- Member companies and individual members may use this material in presentations and literature under the following conditions:
 - Any slide or slides used must be reproduced in their entirety without modification
 - SNIA must be acknowledged as the source of any material used in the body of any document containing material from these presentations.
- This presentation is a project of SNIA.
- Neither the author nor the presenter is an attorney and nothing in this presentation is intended to be, or should be construed as legal advice or an opinion of counsel. If you need legal advice or a legal opinion please contact your attorney.
- The information presented herein represents the author's personal opinion and current understanding
 of the relevant issues involved. The author, the presenter, and the SNIA do not assume any
 responsibility or liability for damages arising out of any reliance on or use of this information.

NO WARRANTIES, EXPRESS OR IMPLIED. USE AT YOUR OWN RISK.



Agenda

- Edge AI Use Cases and Lifecycle
- Edge Data What is New/Different?
- Federated Learning at the Edge
- Dealing with Not Having Shared Storage
- Privacy for Edge AI





Edge Al

Use Cases and Lifecycle



Why AI at the Edge?

• Why?

- Latency
- Privacy
- Bandwidth

Challenges

- Shift from inference only to the full continuum
- Edge applications =
 - Data pipeline (ingest, infer, publish, store, visualize)
 - Business logic (send an alert, stop a process, take an action)





Phases of AI at the Edge



TECHNOLOGIES

CSTI

AI Edge Use Cases



How can I better PREDICT AND Reduce Downtime?



How can I Optimize Operation for higher throughput?



How can I meet rising requirement on product quality?



How can I leverage AI for better business outcomes?



Edge Data

What is New/Different?



Data

Edge data

- Highly distributed
- Crosses boundaries (OT)
- Heterogenous compute
- Dynamic

Cloud data

- Aggregated
- One policy (IT)
- Elastic abstracted compute
- Static (highly available)



Edge Data

What is needed?

- Elastic compute
- Data abstractions (multimodal support, availability, reliability, redundancy, ...)
- Ease of testing and experimentation

How to get there?

- Advances in AI are proving that making sense of data is very valuable (LLMs)
- Rethinking the persona (democratizing AI)
- Al (and data) as a tool



Federated Learning at the Edge

Machine learning without sharing data?



What is Federated Learning?

- Multiple parties
- Train a machine learning model
- Collaboratively
- Without sharing training data



Neonatal Care

 Innocens predicts sepsis incidents

 More data from more hospitals

 Federated gradient boosted trees





Environmental Monitoring

- Visual inspection
- Oil fields are typically shared
- Pipeline of DNNs





Aggregator (A)







Basic Federated Learning 1. Aggregator queries each parties about information Aggregator (A) necessary for learning a predictive model. (e.g. Weights, Gradients, Samples Counts). Q Q Q D_N D_2 D₁ Party 1 (P₁) Party N (P_N) Party 2 (P_2)























Dealing with Not Having Shared Storage



Hyperparameter Tuning for Federated Learning

- Hyperparameters are important for model performance.
- HPO is typically automated or assisted.
- Federated Learning makes
- HPO harder due to lack of data access
- Introduces new HPs



Ref: McMahan, Brendan, et al. "Communication-efficient learning of deep networks from decentralized data." *Artificial Intelligence and Statistics*. PMLR, 2017



Single-shot HPO Framework: FLoRA

Aggregator (A)

HPO rea

 D_2

• One extra round of communication to collect information (local HPO results)

όόό

Party 1

• One complete FL training

3. Send per-party set of (HP, loss) pairs to R₁ aggregator 1. Send the local HPO request to parties

Rn

HPO rea

4. Generate **unified loss surface** using collected info to map any HP to a loss, and select a **single HP** that minimizes this loss surface

5. Execute **single federated training** with selected HP

2. Run independent local HPOs to generate set of (HP, loss) pairs



HPO rea

 R_2

Ref.: Y Zhou, P Ram, T Salonidis, N Baracaldo, H Samulowitz, H Ludwig. "Single-shot general hyper-parameter optimization for federated learning.". In International Conference on Learning Representations, 2023

òòċ

Party p

Fairness and Bias Mitigation in Federated Learning

- Problem:
 - Traditional methods to analyze and mitigate statistical bias in ML require full training dataset
 - Data locations introduce new sources of bias
- Solution:
 - Pre-processing: Local reweighing
 - In-processing: Federated prejudice removal



Privacy for Edge AI



Inference Threats

 In untrusted environments adversaries may try to infer information by analyzing other parties replies





- 2. Gradient of a bag of words: non-zero means the data has a word
- 3. Properties e.g., was someone wearing glasses [5]

Privacy Techniques

- Differential privacy,
- Fully Homomorphic,
- Functional encryption with TPA,
- Functional encryption using decentralized MIFE,
- Threshold Pailler,
- Pailler,
- Secret Sharing and
- **Trusted Execution Environments**





Combine multi-party computation and differential privacy to achieve higher model accuracy through reduced noise





Multi-Party Computation

• Provides privacy of the inputs



Depending on the crypto system, the resulting aggregated model may be in plaintext or ciphertext

Combination Approach Outperforms Baselines



- The more parties the data is divided amongst, the more noise the *Local Differential Privacy* approach requires, resulting in poor accuracy
- In contrast, our approach produces stable accuracy
- The baseline converges with the "random guess" line at around 100 parties

Results shown for decision tree, Nursery dataset from UCI, $\epsilon = 0.5$, 10 parties

© C34 C SNIA, All Rights Reserved.

A Hybrid Approach to Privacy-Preserving Federated Learning Truex et al (Best paper award)









- Al is changing edge compute forever; edge is also changing Al
- Data at the edge is complex, noisy, and distributed
- Solving the data challenges is the key to democratizing AI

Join us in innovating in edge Al!









Thanks for Viewing this Webinar

- Please rate this presentation and provide us with feedback
- This webinar and a copy of the slides are available at the SNIA Educational Library <u>https://www.snia.org/educational-library</u>
- A Q&A from this webinar will be posted to the SNIA Cloud blog: <u>www.sniacloud.com/</u>
- Follow us on X (formerly) Twitter <u>@SNIACloud</u>



Thank You!

