

## **Temporal Predictive Coding for Gradient Compression in Distributed Learning**

Adrian Edin<sup>1</sup>, Zheng Chen<sup>1</sup>, Michel Kieffer<sup>2</sup> and Mikael Johansson<sup>3</sup>

<sup>1</sup>Div. of Communication Systems, Dept. of Electrical Engineering (ISY), Linköping University (LiU), Sweden <sup>2</sup>Université Paris-Saclay, CentraleSupelec, CNRS, Laboratoire des Signaux et Systèmes, Gif-sur-Yvette, France <sup>3</sup>School of Electrical Engineering and Computer Science, KTH, Stockholm, Sweden





## Conclusions

 $10^{5}$ 

 $10^{6}$ 

# bits

 $[\times 10^5]$ 

135.20

12.49

8.60

6.54

8.40

13.00

Total transmitted bits

# iters.

731

691

726

733

733

733

We proposed a prediction-based gradient compression with event-triggered communication for distributed learning. Remarks:

- · The current gradient can be actively predicted using previous gradient information.
- Omitting transmission when prediction is accurate significantly reduce the amount of transmitted bits.

 $oldsymbol{e}_k^{(t)} = oldsymbol{g}_k^{(t)} - \widehat{oldsymbol{g}}_k^{(t)}$ 



**Prediction residual** 

Properties from the LS-estimate:

• Orthogonality  $\widehat{g}_{k}^{(t)} \perp e_{k}^{(t)}$  • Norm reduction:  $\left\| e_{k}^{(t)} \right\| \leq \left\| g_{k}^{(t)} \right\|$ 

Block Diagram for Agent k (Each agent is processed independently)

- Synchronized memory updates use imperfect gradient  $\widetilde{g}_{k}^{(t)}$
- Prediction coefficients  $a^{*(t)}$ transmitted with high bit count, *i.e.* negligible distortion.
- Both agent and PS knows  $\widehat{g}_{L}^{(t)}$

This work was presented at the 60th Allerton conference on Communication, Control, and Computing, Illinois, USA, 2024.

Link to the paper and more information: https://ostpopcorn.se/allerton24

