

# Streaming Tensor Factorization for Infinite Data Sources

Shaden Smith - Intel Parallel Computing Lab

Kejun Huang - University of Minnesota

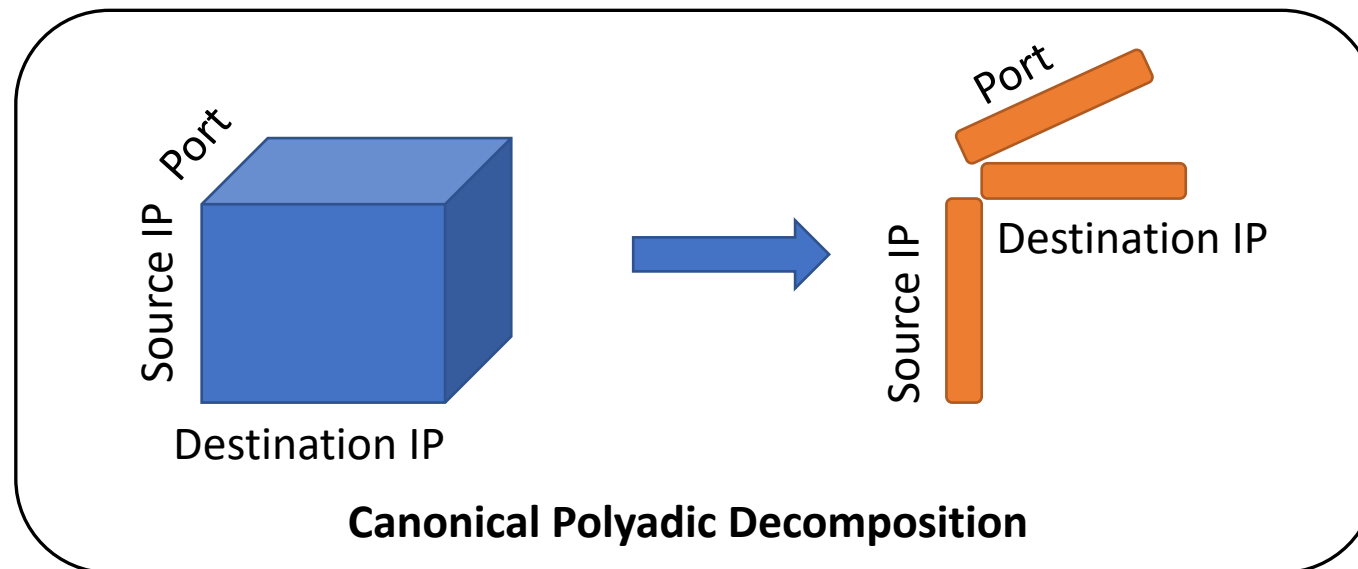
Nicholas D. Sidiropoulos - University of Virginia

George Karypis - University of Minnesota

`Shaden.Smith@intel.com`

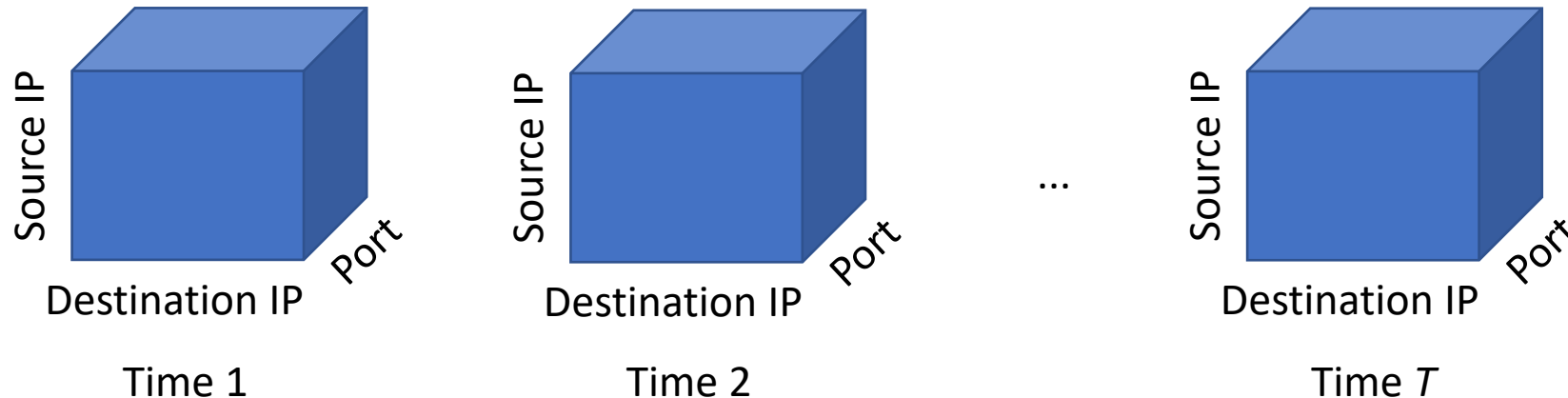
# Tensor factorization

- Multi-way data can be naturally represented as a *tensor*.
- Tensor factorizations are powerful tools for facilitating the analysis of multi-way data.
  - Think: singular value decomposition, principal component analysis.



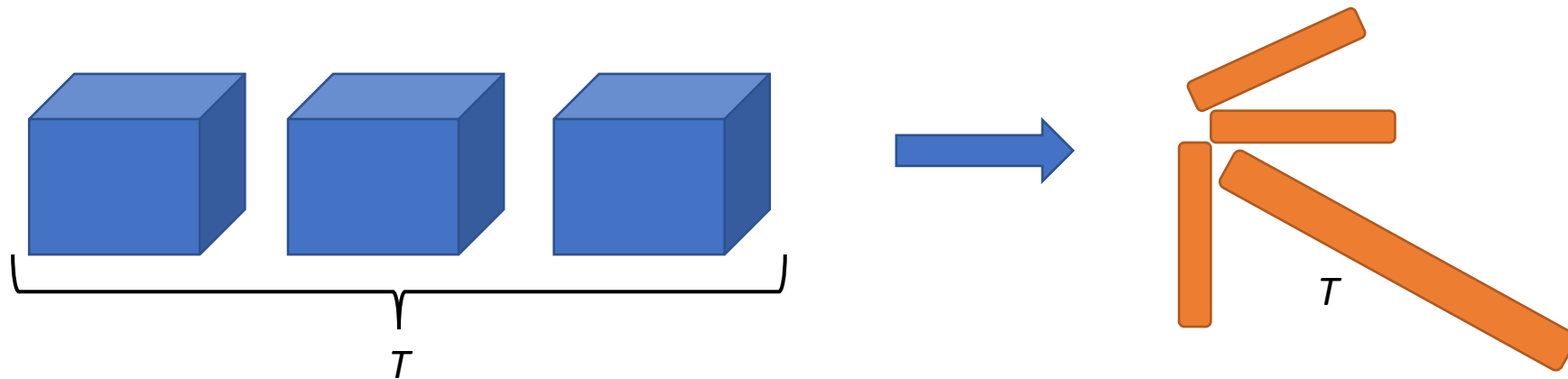
# Streaming data

- We often need to analyze multi-way data that is *streamed*.
  - Applications include: cybersecurity, discussion tracking, traffic analysis, video monitoring, ...
- A batch of data arrives each timestep  $1, \dots, T$ .
  - **$T$  may be infinite!**
- Batches are assumed to come from the same generative model.
  - In practice, we must account for the model slowly changing over time.



# Streaming tensor factorization

- The collection of  $N$ -dimensional tensors can be viewed as an  $(N+1)$ -dimensional tensor observed over time.
- We want to cheaply *update* an existing factorization each timestep to incorporate the latest batch of data.
- **Challenge:** storing historical tensor *or* factorization data that grows with time is infeasible.
- **Challenge:** we would like to apply constraints such as non-negativity to the factorization.



# CP-stream: optimization problem

- We start from the following non-convex optimization problem over all timesteps:

$$\begin{aligned} & \underset{\{\mathbf{A}^{(n)} \in \mathbb{R}^{I_n \times K}\}, \{\mathbf{s}_t \in \mathbb{R}^K\}}{\text{minimize}} && \frac{1}{2} \sum_{t=1}^T \left( \left\| \mathbf{x}_t - \llbracket \{\mathbf{A}^{(n)}\}; \mathbf{s}_t \rrbracket \right\|^2 + \lambda \|\mathbf{s}_t\|^2 \right) \\ & \text{subject to} && \mathbf{A}^{(n)} \in \mathcal{C}, \end{aligned}$$

- We constrain the factor matrices to have column norms  $\leq 1$ .
  - This improves stability due to a scaling ambiguity in the CPD.
- The  $\mathbf{s}_t \in \mathbb{R}^K$  vectors form the rows of  $\mathbf{S}$ , the temporal factor matrix.

# CP-stream: formulation

- To avoid storing historic tensor data, we follow (Vandecappelle et al. 2017) and instead use the historical factorization:

$$\begin{aligned} & \underset{\{\mathbf{A}^{(n)}\}, \mathbf{s}_t}{\text{minimize}} && \frac{1}{2} \left\| \mathcal{X}_t - \llbracket \{\mathbf{A}^{(n)}\}; \mathbf{s}_t \rrbracket \right\|^2 + \sum_{i=1}^{t-1} \frac{\mu^{t-i}}{2} \left\| \llbracket \{\mathbf{A}_{t-1}^{(n)}\}; \mathbf{s}_i \rrbracket - \llbracket \{\mathbf{A}^{(n)}\}; \mathbf{s}_t \rrbracket \right\|^2 \\ & \text{subject to} && \mathbf{A}^{(n)} \in \mathcal{C}, \end{aligned}$$

- $\mu$  is a *forgetting factor* used to down-weight the importance of older data.
- **Limitation:** this still requires  $\mathbf{S} \in \mathbb{R}^T \times K$ .

# CP-stream: algorithm (details in paper/poster)

When a new batch of data arrives at time  $t$ :

1. Compute  $\mathbf{s}_t$ . This has a closed-form solution involving the new batch of tensor data and the previous factor matrices.
  - Complexity does not depend on  $T$ .
2. Update the factor matrices. We use alternating optimization with ADMM (AO-ADMM; Huang & Sidiropoulos 2016).
  - The temporal factor  $\mathbf{S}$  is only used in its **compact** Gramian form  $\mathbf{S}^T \mathbf{S}$ , which is computed recursively:

$$\mathbf{G}_t = \mu \mathbf{G}_{t-1} + \mathbf{s}_t \mathbf{s}_t^T.$$

# Extensions

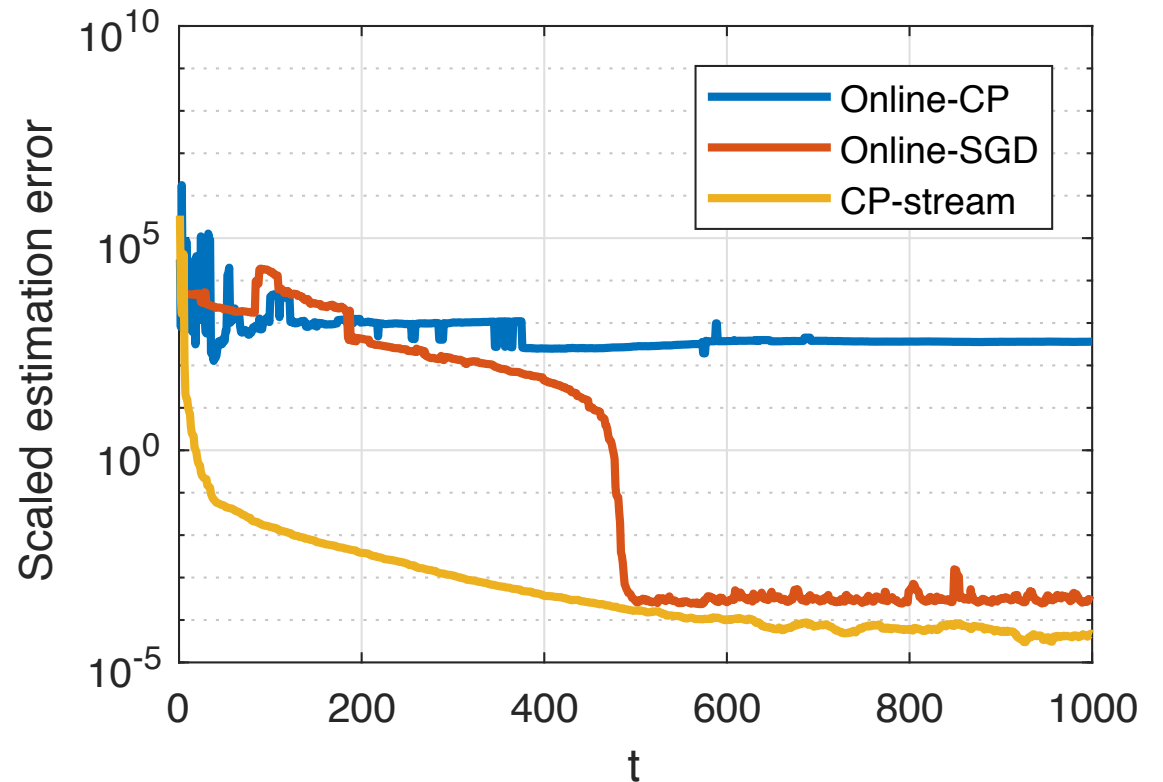
- CP-stream supports additional constraints/regularizations. For stability, they are combined with the column norm constraint (*proof of convergence in paper*).
  - Non-negativity
  - $\ell_1$  regularization to promote sparse factors
- Tensor sparsity:
  - CP-stream scales linearly in the number of non-zeros and makes use of the existing optimized kernels.
  - Sparsity is not treated as *missing*, because absence of activity also carries meaning in our applications.



# Evaluation

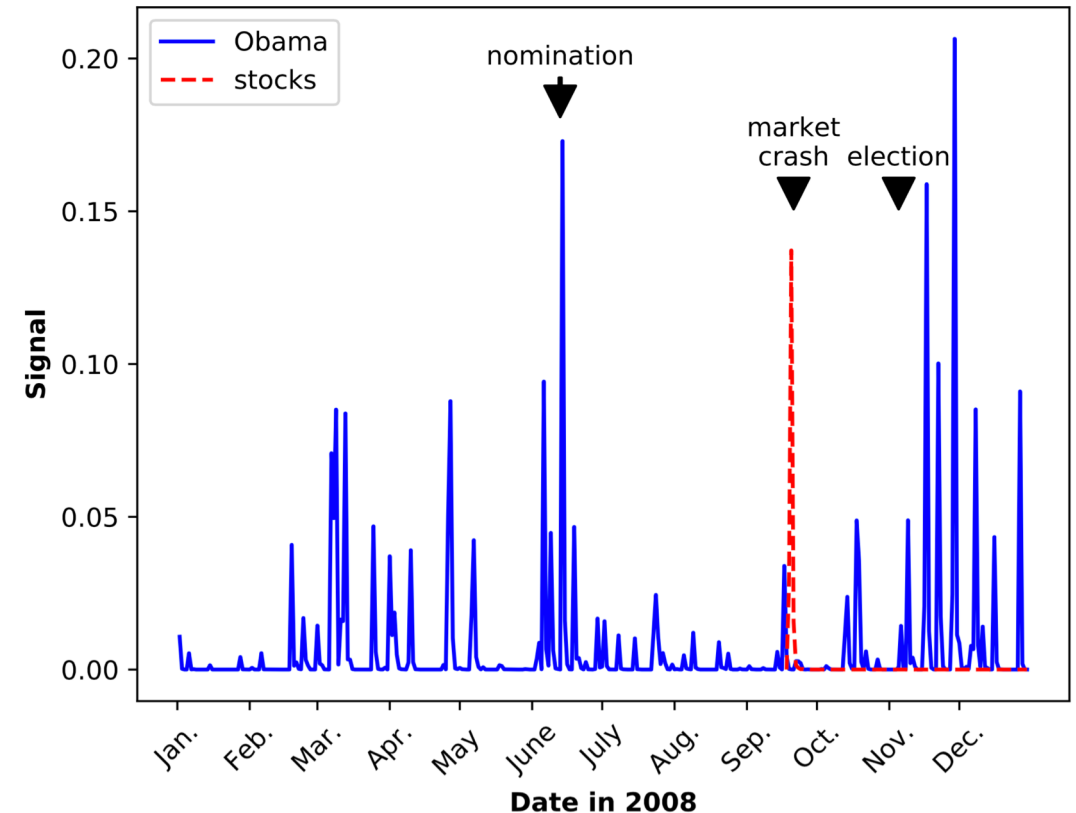
- We generated a dense 100x100x1000 tensor from rank-10 factors (plus noise).
- We compare against:
  - Online-CP (Zhou et al., 2016)
  - Online-SGD (Mardani et al., 2015)
- Shown is the estimation error of the known ground-truth factors:

$$\frac{\|\mathbf{A}_{\natural}^{(1)} - \mathbf{A}_t^{(1)}\|^2}{\|\mathbf{A}_{\natural}^{(1)}\|^2} + \frac{\|\mathbf{A}_{\natural}^{(2)} - \mathbf{A}_t^{(2)}\|^2}{\|\mathbf{A}_{\natural}^{(2)}\|^2}$$



# Case study: discussion tracking

- Comments on reddit.com form a *(user, community, word)* tensor.
  - A new batch arrives each day.
  - 65M non-zeros over one year.
- Each user, community, and word are represented by a low-rank vector in the factorization.
- Tracking the vectors representing the word “Obama” and the stocks community reveals events in 2008.



# Wrapping up

- Streaming tensor factorization has applications in areas such as cybersecurity, discussion tracking, and traffic analysis.
- CP-stream uses a formulation suitable for long-term streaming, and supports sparsity and constraints.
- Our source code is to be open sourced as part of SPLATT
  - <https://github.com/ShadenSmith/splatt>
- Sparse tensor datasets available in FROSTT:
  - <http://frostt.io/>
- **Contact:** [Shaden.Smith@intel.com](mailto:Shaden.Smith@intel.com) **or** `shaden@cs.umn.edu`

Backup

# AO-ADMM

$$(4.3) \quad \mathbf{s}_t \leftarrow \left( \bigotimes_{n=1}^N \mathbf{A}_{t-1}^{(n)\top} \mathbf{A}_{t-1}^{(n)} + \lambda \mathbf{I} \right)^{-1} \left( \bigodot_{n=1}^N \mathbf{A}_{t-1}^{(n)} \right)^\top \text{vec}(\mathcal{X}_t).$$

$$(4.5) \quad \Phi^{(n)} = \left( \bigotimes_{\nu \neq n} \mathbf{A}^{(\nu)\top} \mathbf{A}^{(\nu)} \right) \circledast (\mu \mathbf{G}_{t-1} + \mathbf{s}_t \mathbf{s}_t^\top),$$

$$(4.6) \quad \Psi^{(n)} = \left( \bigodot_{\nu \neq n} \mathbf{A}^{(\nu)} \right)^\top \text{vec}(\mathcal{X}_t) + \mathbf{A}^{(n)} \left( \left( \bigotimes_{\nu \neq n} \mathbf{A}_{t-1}^{(\nu)\top} \mathbf{A}_{t-1}^{(\nu)} \right) \circledast \mu \mathbf{G}_{t-1} \right),$$

and

$$\mathbf{G}_{t-1} = \sum_{i=1}^{t-1} \mu^{t-i} \mathbf{s}_i \mathbf{s}_i^\top.$$

---

## Algorithm 1 CP-stream

---

**Require:**  $\mathcal{X}_1, \mathcal{X}_2, \dots, \mathcal{X}_T$ ; forgetting factor  $\mu$

- 1: initialize  $\mathbf{A}_0^{(1)}, \dots, \mathbf{A}_0^{(N)}$
  - 2:  $\mathbf{G}_0 \leftarrow \mathbf{0}$
  - 3: **for**  $t = 1, \dots, T$  **do**
  - 4:      $\mathbf{s}_t \leftarrow$  least-squares update (4.3)
  - 5:     **repeat**
  - 6:         **for**  $n = 1, \dots, N$  **do**
  - 7:             construct  $\Phi^{(n)}$  and  $\Psi^{(n)}$  per (4.5) and (4.6)
  - 8:              $\rho = \text{tr}(\Phi^{(n)}) / K$
  - 9:              $\mathbf{A}_t^{(n)} \leftarrow$  ADMM iterates (4.7)
  - 10:         **end for**
  - 11:     **until** convergence
  - 12:      $\mathbf{G}_t = \mu \mathbf{G}_{t-1} + \mathbf{s}_t \mathbf{s}_t^\top$
  - 13: **end for**
-

# AO-ADMM (2)

$$(4.7) \quad \begin{cases} \tilde{\mathbf{A}} \leftarrow \left( \boldsymbol{\Psi}^{(n)} + \rho(\mathbf{A}^{(n)} + \mathbf{U}) \right) \left( \boldsymbol{\Phi}^{(n)} + \rho \mathbf{I} \right)^{-1}, \\ \mathbf{A}^{(n)} \leftarrow \text{Proj}_{\mathcal{C}} \left[ \tilde{\mathbf{A}} - \mathbf{U} \right], \\ \mathbf{U} \leftarrow \mathbf{U} + \tilde{\mathbf{A}} - \mathbf{A}^{(n)}, \end{cases}$$

---

## Algorithm 1 CP-stream

---

**Require:**  $\mathcal{X}_1, \mathcal{X}_2, \dots, \mathcal{X}_T$ ; forgetting factor  $\mu$

- 1: initialize  $\mathbf{A}_0^{(1)}, \dots, \mathbf{A}_0^{(N)}$
  - 2:  $\mathbf{G}_0 \leftarrow \mathbf{0}$
  - 3: **for**  $t = 1, \dots, T$  **do**
  - 4:      $\mathbf{s}_t \leftarrow$  least-squares update (4.3)
  - 5:     **repeat**
  - 6:         **for**  $n = 1, \dots, N$  **do**
  - 7:             construct  $\boldsymbol{\Phi}^{(n)}$  and  $\boldsymbol{\Psi}^{(n)}$  per (4.5) and (4.6)
  - 8:              $\rho = \text{tr}(\boldsymbol{\Phi}^{(n)})/K$
  - 9:              $\mathbf{A}_t^{(n)} \leftarrow$  ADMM iterates (4.7)
  - 10:         **end for**
  - 11:     **until** convergence
  - 12:      $\mathbf{G}_t = \mu \mathbf{G}_{t-1} + \mathbf{s}_t \mathbf{s}_t^\top$
  - 13: **end for**
-