

---

---

CSCI-6961 Project G3

# **Asynchronous Distributed ADMM for Consensus Optimization**

— Md Shamim Hussain, Joshua Kavner, Osama  
Minhas, Allie Solko-Breslin, Lingyu Zhang —

**Dec 2, 2020**

---

---

# Talk Outline

1. Motivation
2. Synchronous (Centralized) ADMM
3. Asynchronous (Centralized) ADMM
4. Network Topology Analysis
5. Empirical Evaluation
6. Discussion



# Motivation

## ADMM vs Mini-batch SGD

- Jointly optimizes data on multiple servers without exchanging data
- Consensus optimization with ADMM instead of Mini-batch SGD
- Parallelization of optimization using independent different workers

## Distributed Asynchronous ADMM

- Avoids straggling by allowing workers to update at different speeds
- Bounds delay on slowest worker to encourage regular updating

# Related Works

ADMM -- Gabay et al 1976



126 page survey by  
Boyd et. al. 2011

$O(1/t)$  convergence ADMM  
[Bingsheng He, Xiaoming Yuan, SIAM 2011]

Decentralized ADMM  
[Ermin Wei, Asuman Ozgander, IEEE 2012]

Asynchronous Decentralized ADMM  
[Ermin Wei, Asuman Ozgander, IEEE 2013]

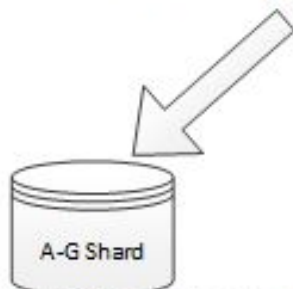
**Async. Centralized ADMM**  
[Ruiliang Zhang, James Kwok, ICML 2014]

# Main result of the paper

- Asynchronous distributed processing can be extended to ADMM
- Faster than synchronous ADMM in practice
- Models that can be formulated using ADMM can be learned efficiently
- $O(1/t)$  convergence

# ADMM: Alternating Direction Method of Multipliers

Key	Name	Description	Stock	Price	LastOrdered
ARC1	Arc welder	250 Amps	8	119.00	25-Nov-2013
BRK8	Bracket	250mm	46	5.66	18-Nov-2013
BRK9	Bracket	400mm	82	6.98	1-Jul-2013
HOS8	Hose	1/2"	27	27.50	18-Aug-2013
WGT4	Widget	Green	16	13.99	3-Feb-2013
WGT6	Widget	Purple	76	13.99	31-Mar-2013

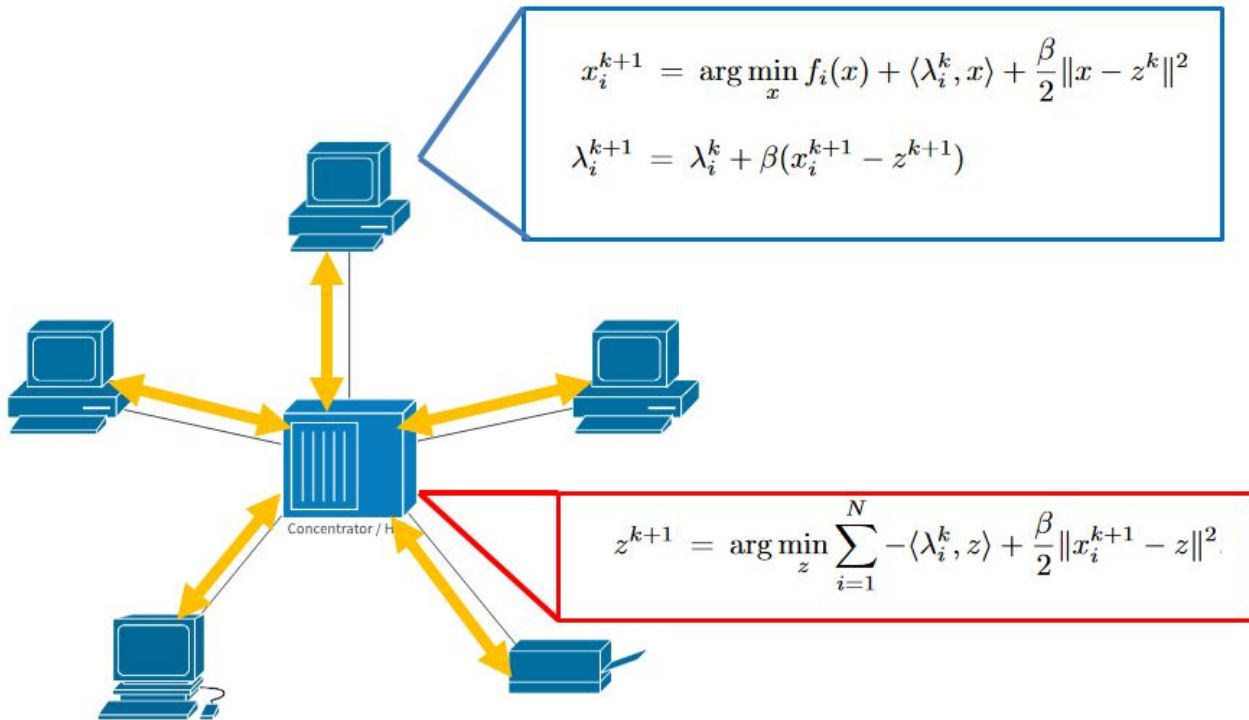


Key	Name	Description	Stock	Price	LastOrdered
ARC1	Arc welder	250 Amps	8	119.00	25-Nov-2013
BRK8	Bracket	250mm	46	5.66	18-Nov-2013
BRK9	Bracket	400mm	82	6.98	1-Jul-2013



Key	Name	Description	Stock	Price	LastOrdered
HOS8	Hose	1/2"	27	27.50	18-Aug-2013
WGT4	Widget	Green	16	13.99	3-Feb-2013
WGT6	Widget	Purple	76	13.99	31-Mar-2013

# Synchronous (Centralized) ADMM



# ADMM: Alternating Direction Method of Multipliers

Suppose we have the following problem:

$$\text{minimize } f(x) + g(z) \text{ such that } Ax + Bz = c.$$

Augmented Lagrangian



The augmented Lagrangian is of the form:

$$L_{\beta}(x, z, \lambda) = \underbrace{f(x)}_{\text{red}} + \underbrace{g(z)}_{\text{blue}} + \underbrace{\lambda^T (Ax + Bz - c) + (\beta/2) \| (Ax + Bz - c) \|_2^2}_{\text{green}}.$$

Iterations



ADMM consists of the iterations:

$$x^{k+1} = \operatorname{argmin}_x L_{\beta}(x, z^k, \lambda^k)$$

$$z^{k+1} = \operatorname{argmin}_z L_{\beta}(x^k, z, \lambda^k)$$

$$\lambda^{k+1} = \lambda^k + \beta(Ax^{k+1} + Bz^{k+1} - c)$$



# ADMM: Consensus Optimization

Let  $(x_i)_{i=1}^N \in \mathbb{R}^n$ . Then  $x = [x_1 || \dots || x_N]^T \in \mathbb{R}^{nN}$

Substitute  $A \leftarrow \text{diag}_N(\mathbb{I}_{n \times n}) = \begin{pmatrix} \mathbb{I} & 0 & \dots & 0 \\ 0 & \mathbb{I} & \dots & 0 \\ \dots & \dots & \dots & \dots \\ 0 & 0 & \dots & \mathbb{I} \end{pmatrix}$

$B \leftarrow -1 * \text{repeat}_N(\mathbb{I}_{n \times n}) = -(\mathbb{I} \quad \mathbb{I} \quad \dots \quad \mathbb{I})^T$

$c \leftarrow 0$

$f(x) := \sum_{i=1}^N f_i(x_i)$

$g(z) := 0$

minimize  $f(x) + g(z)$   
such that  $Ax + Bz = c$

# ADMM: Consensus Optimization



$$\min_{x_1, \dots, x_N, z} \sum_{i=1}^N f_i(x_i) : x_i = z, i = 1, 2, \dots, N$$

Augmented Lagrangian



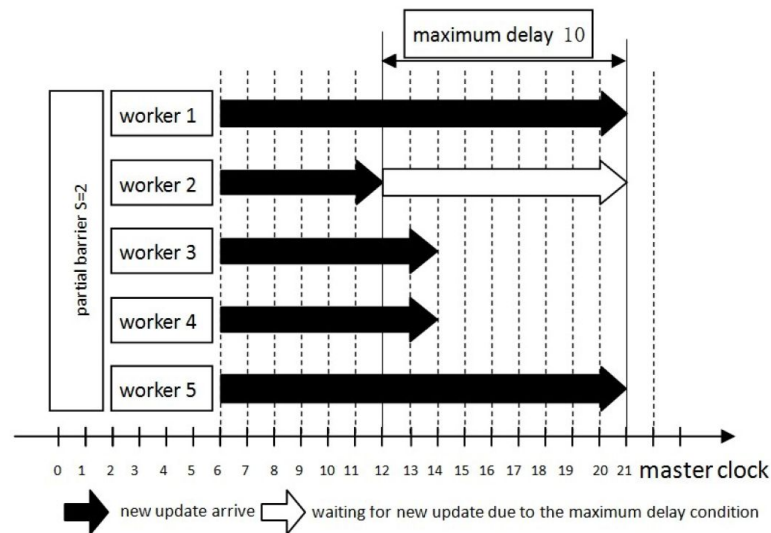
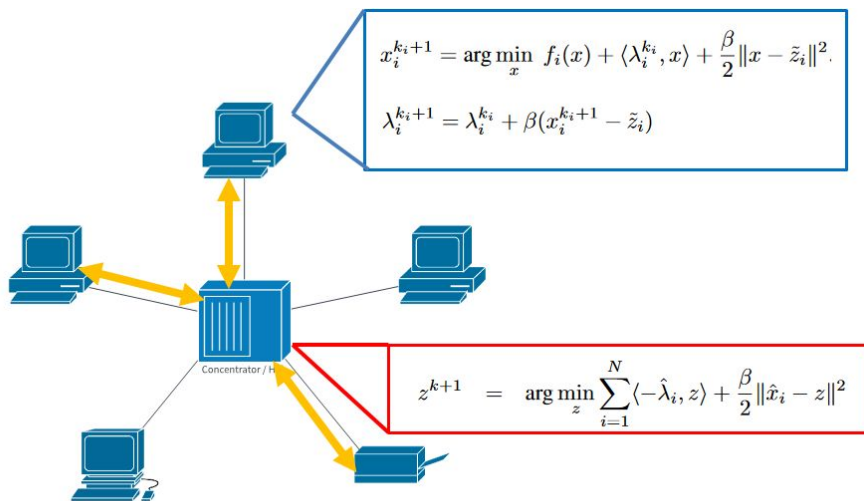
$$L(\{x_i\}, z) = \sum_{i=1}^N f_i(x_i) + \langle \lambda_i, x_i - z \rangle + \frac{\beta}{2} \|x_i - z\|^2$$

Iterations



$$\begin{aligned} x_i^{k+1} &= \arg \min_x f_i(x) + \langle \lambda_i^k, x \rangle + \frac{\beta}{2} \|x - z^k\|^2 \\ z^{k+1} &= \arg \min_z \sum_{i=1}^N -\langle \lambda_i^k, z \rangle + \frac{\beta}{2} \|x_i^{k+1} - z\|^2 \\ \lambda_i^{k+1} &= \lambda_i^k + \beta(x_i^{k+1} - z^{k+1}). \end{aligned}$$

# Asynchronous (Centralized) ADMM



**Partial Barrier and Bounded Delay**

# Asynchronous (Centralized) ADMM: Algorithm

## Master

Repeat:

- Wait for
  - A minimum of  $S$  updates from different workers
  - No update delayed more than  $\tau$  steps
- Update consensus variable  $z$
- Broadcast consensus variable  $z$  to workers

## Worker $i$

Repeat:

- Update local variable  $x_i$
- Send local variable  $x_i$  and multiplier  $\lambda_i$  to master
- Wait for updated consensus variable  $z$
- Update multiplier  $\lambda_i$

# Asynchronous (Centralized) ADMM: Analysis

**Theorem 4.2** Let  $(x^*, z^*)$  be the optimal (primal) solutions of the minimization problem (2), and  $\{\lambda_i^*\}_{i=1}^N$  the corresponding optimal dual solution. Consider  $\bar{x}_i = \frac{1}{T_i} \sum_{t=0}^{T_i-1} x_i^t$  and  $\bar{z} = \frac{1}{T} \sum_{t=0}^{T-1} z^t$ . Then

$$\mathbb{E} \left[ \sum_{i=1}^N \underbrace{f_i(\bar{x}_i) - f_i(x^*)}_{\text{blue}} + \underbrace{\langle \lambda_i^*, \bar{x}_i - \bar{z} \rangle}_{\text{green}} \right] \leq$$
$$\frac{N\tau}{2TS} \left[ \sum_{i=1}^N \beta \|z_i^0 - z^*\|^2 + \frac{1}{\beta} \|\lambda_i^0 - \lambda_i^*\|^2 \right]$$

where  $z_i^0$  and  $\lambda_i^0$  are the initial values of  $z_i$  and  $\lambda_i$  respectively, at worker  $i$ .

# Comparative Performance Guarantees

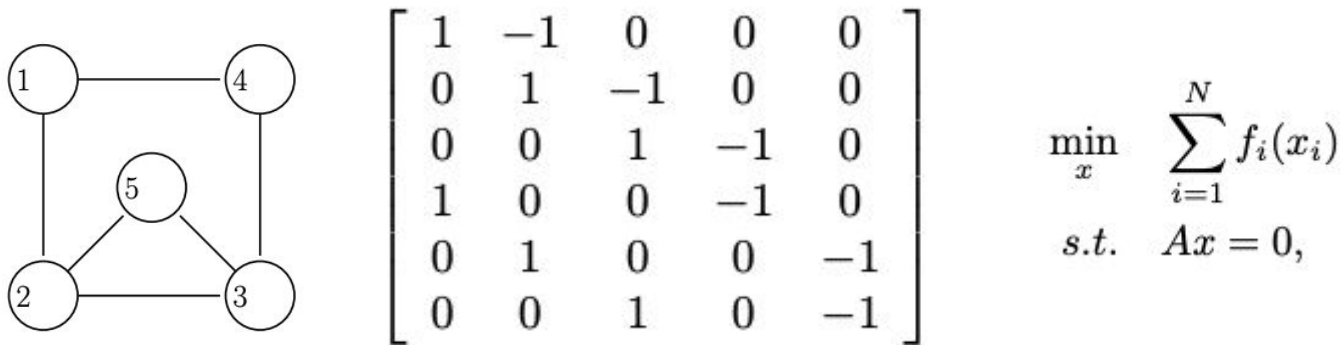
	<b>Synchronous</b>	<b>Asynchronous</b>
<b>Centralized</b>	$O(1/t)$ [He, Yuan, 2011]	$O(1/t)$ [Zhang, Kwok, 2014]
<b>Decentralized</b>	$O(1/t)$ [Wei, Ozdaglar, 2012]	$O(1/t)$ [Wei, Ozdaglar, 2013]

# From centralized to decentralized ADMM

- Centralized ADMM requires a master node.
  - High computation cost
  - Risk of losing data if master node fails
  - Crowding of messages
- Extending standard ADMM to a decentralized setting:
  - Lower computation cost.
  - Data privacy.
- How do different topologies affect model performance?

# Represent network topology as a matrix

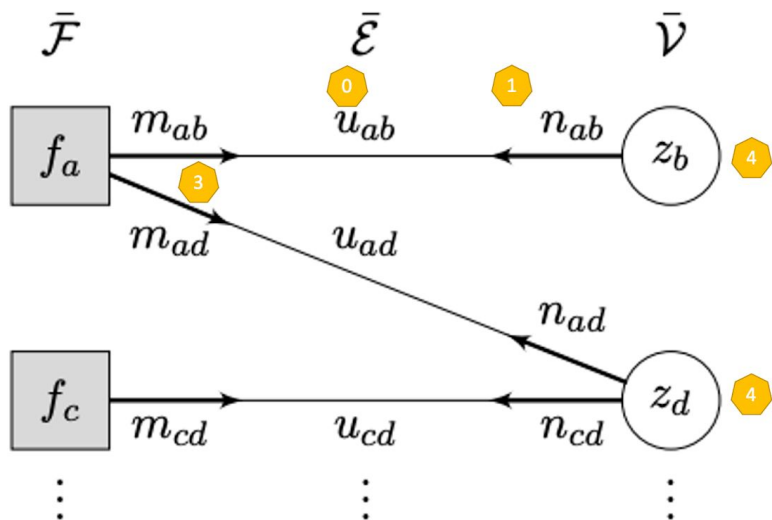
- Represent network topology as a matrix [1].
  - N nodes and M edges.
  - Each node represents an agent.
  - Each node has an associated cost function.
  - Final optimization: To minimize the sum of the loss functions from the N nodes.
  - Incidence matrix is used as a constraint in the optimization task.





# Formulate ADMM as decentralized version.

- Message passing task [1].

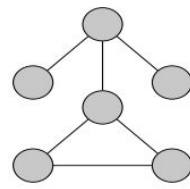
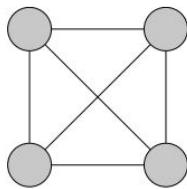
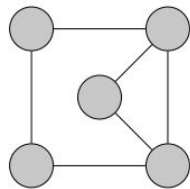
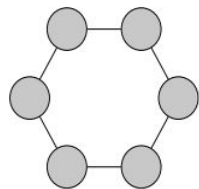


$$\begin{aligned}
 & \text{0} \quad u_{ab}^{t+1} \leftarrow u_{ab}^t + \gamma x_{ab}^{t+1} - z_b^t + (1 - \gamma) z_b^{t-1} && \text{for all } (a, b) \in \bar{\mathcal{E}} \\
 & \text{1} \quad n_{ab}^t \leftarrow z_b^t - u_{ab}^t && \text{for all } (a, b) \in \bar{\mathcal{E}} \\
 & \text{2} \quad x_a^{t+1} \leftarrow \arg \min_{x_a} \left\{ f_a(x_a) + \frac{\rho}{2} \sum_{b \in N_a} (x_{ab} - n_{ab}^t)^2 \right\} && \text{for all } a \in \bar{\mathcal{F}} \\
 & \text{3} \quad m_{ab}^{t+1} \leftarrow \gamma x_{ab}^{t+1} + u_{ab}^t && \text{for all } (a, b) \in \bar{\mathcal{E}} \\
 & \text{4} \quad z_b^{t+1} \leftarrow (1 - \gamma) z_b^t + \frac{1}{|N_b|} \sum_{a \in N_b} m_{ab}^{t+1} && \text{for all } b \in \bar{\mathcal{V}} \\
 & \text{0} \quad u_{ab}^{t+1} \leftarrow u_{ab}^t + \gamma x_{ab}^{t+1} - z_b^{t+1} + (1 - \gamma) z_b^t && \text{for all } (a, b) \in \bar{\mathcal{E}}
 \end{aligned}$$

[1] França, Guilherme, and José Bento. "How is distributed ADMM affected by network topology?." arXiv preprint arXiv:1710.00889 (2017).

# Main factors in network topologies

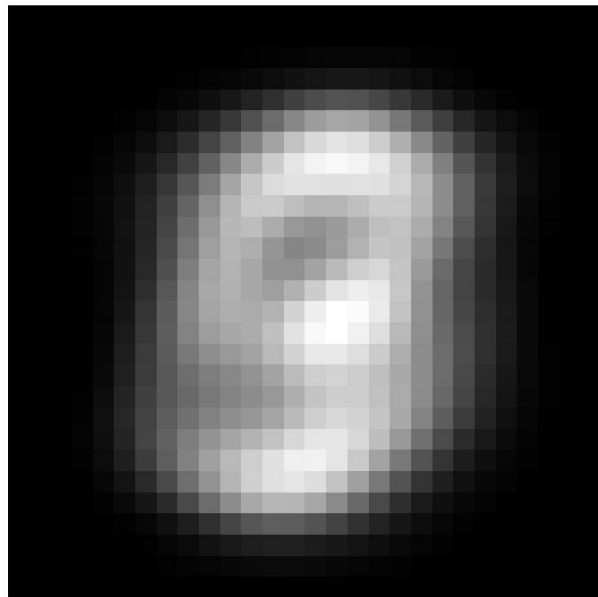
- Main factors in network topologies:
  - Graph has at least one cycle of even length.
  - Graph has a cycle, but not with an even length.
  - Graphs without cycles.
  - Example: periodic grid, ring, k-hop lattice, graph sampled from the Erdos-Renyi model.



# Empirical Evaluation: Data

- Implemented in Python
- TCP communication
- **N=16**, 16 Workers and a Master, each a physical process
- Approximate average image (pixel-by-pixel) over sample from MNIST
- Objective

$$\min_x f(x) = \sum_{i=1}^N \|x - \theta_i\|^2$$



# Empirical Evaluation: Partial Barrier

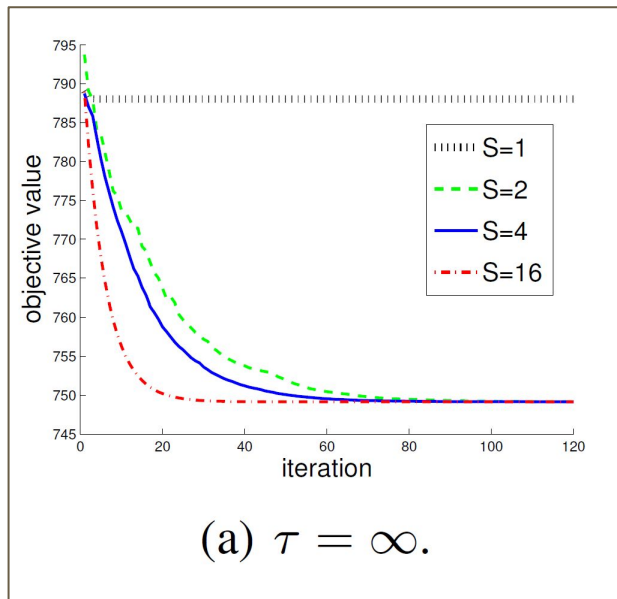
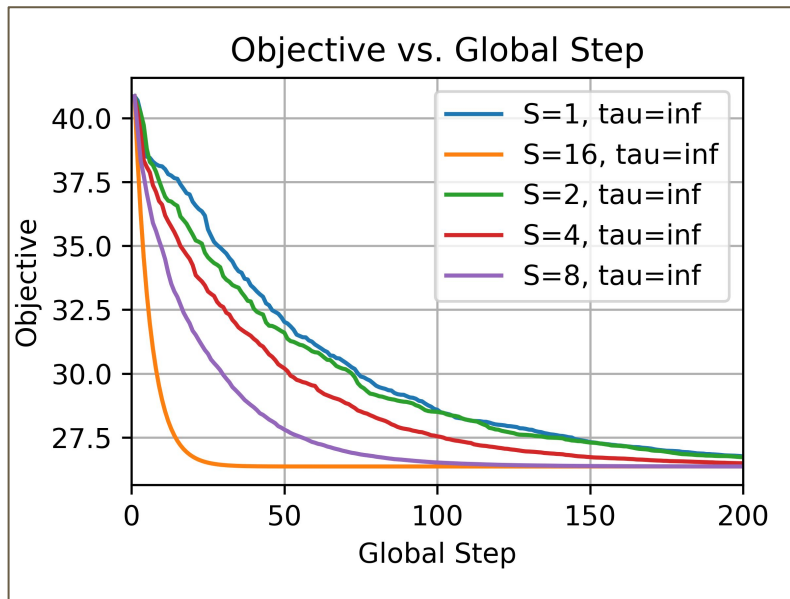


Figure 2a  
Avg loss: random numbers  
[Zhang and Kwok 14]



200 empirical iterations  
Avg loss: MINST sample

# Empirical Evaluation: Partial Barrier

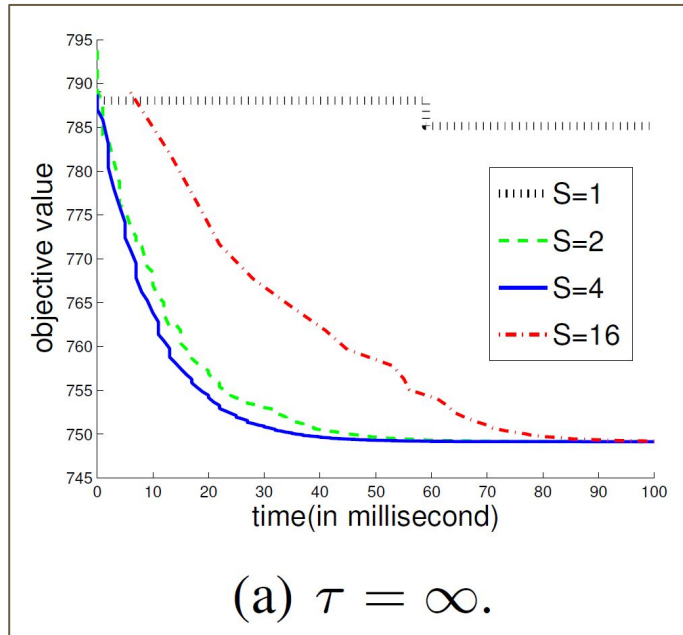
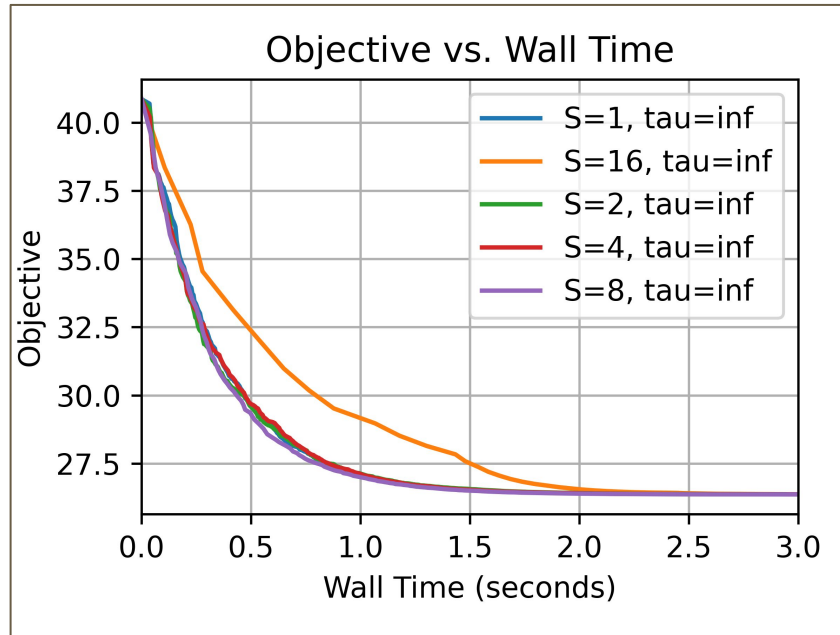


Figure 3a  
Avg loss: random numbers  
[Zhang and Kwok 14]



200 empirical iterations  
Avg loss: MINST sample

# Empirical Evaluation: Bounded Delay

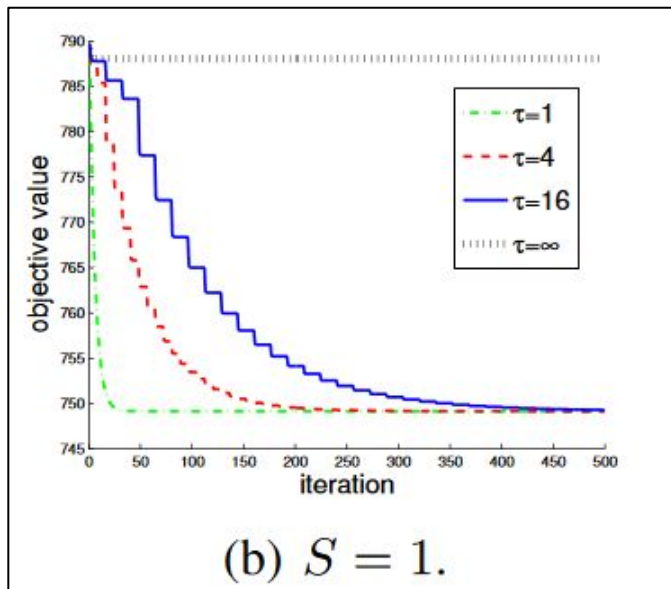
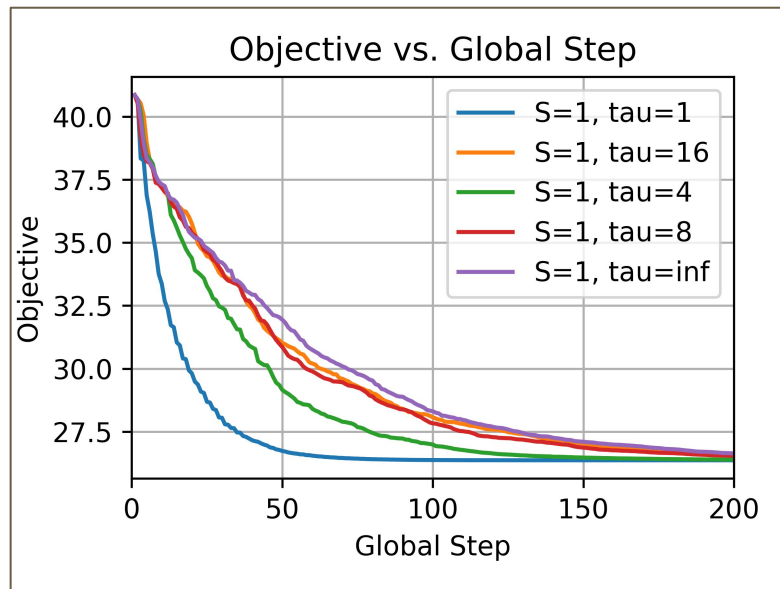


Figure 2b

Avg loss: random numbers  
[Zhang and Kwok 14]



200 empirical iterations

Avg loss: MINST sample

# Empirical Evaluation: Bounded Delay

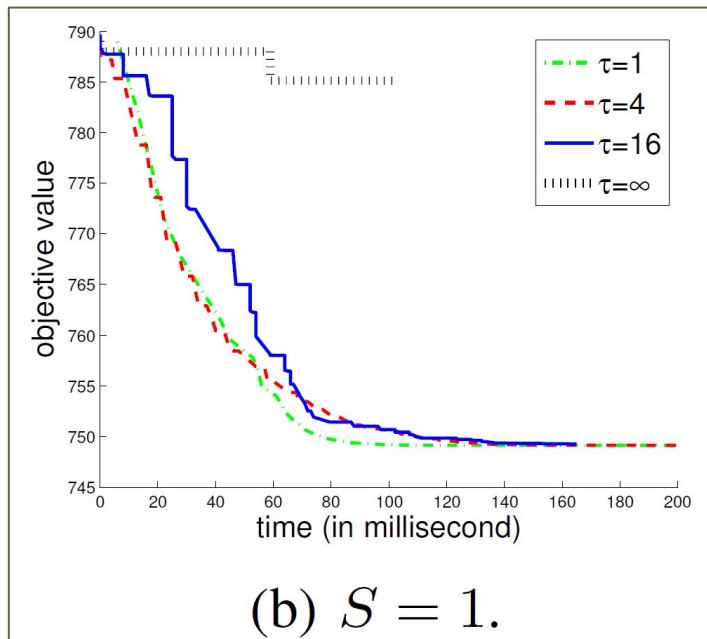
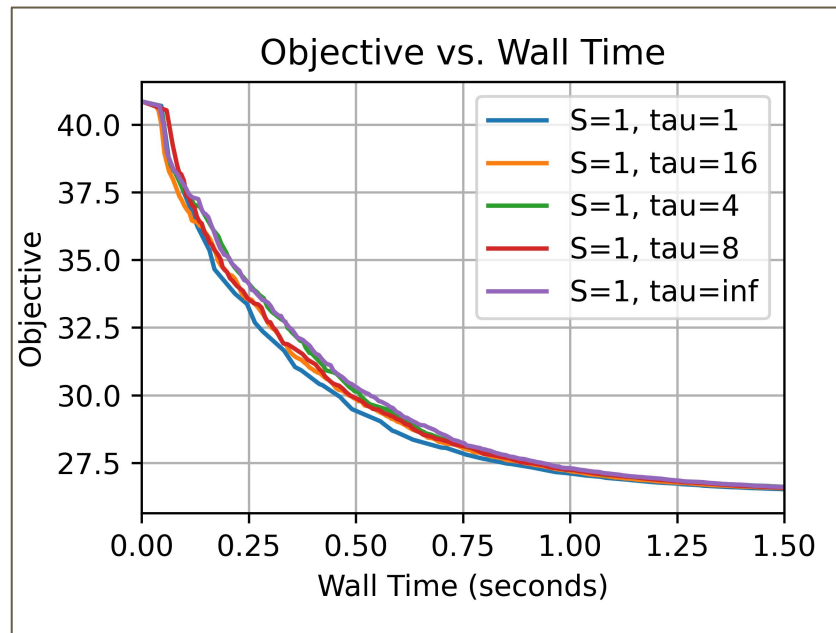


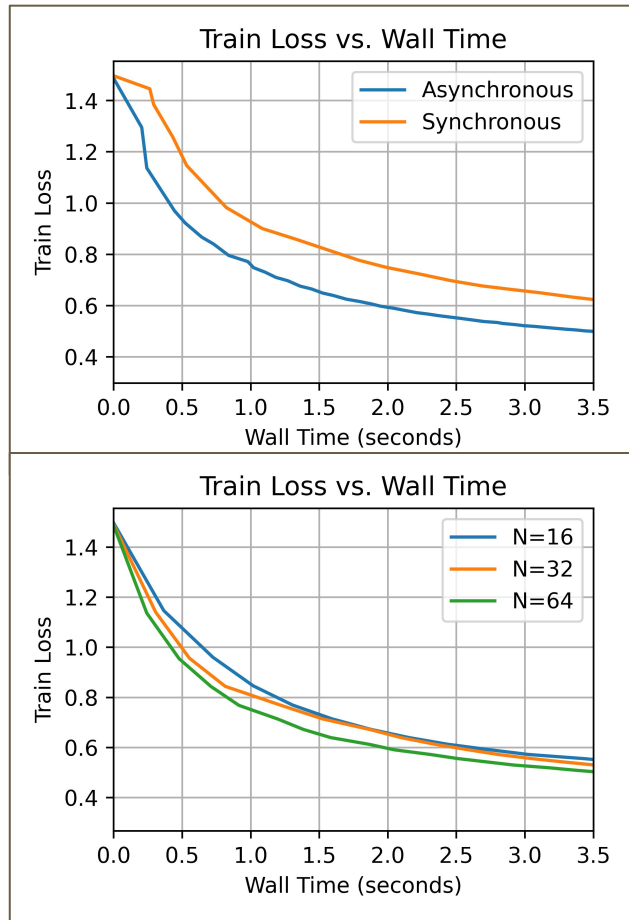
Figure 3b  
Avg loss: random numbers  
[Zhang and Kwok 14]



200 empirical iterations  
Avg loss: MINST sample

# Multiclass Logistic Regression - MNIST

- Each worker does gradient descent to minimize local loss
- Weights and biases as consensus variable
- Faster convergence than synchronous case
- Faster convergence when number of workers is increased







# Advantages

- ADMM's flexible framework allows it to optimize a variety of objective functions and use numerous network topologies.
- By using nodes and distributed memory, ADMM has faster speed and lower computational complexity than other algorithms when working with large data sets.

# Disadvantages

- Synchronous ADMM suffers from straggling - the algorithm is only as fast as its slowest component.
  - Asynchronous ADMM avoids this problem by avoiding the need to wait for the slower components.
- Centralized ADMM suffers from bottlenecking - all nodes communicate with the central node, which can get congested and slow the performance.
  - Decentralized ADMM avoids this problem by having nodes communicate with neighboring nodes instead of the central node.

**Thank You**