

Microsoft Research

Each year Microsoft Research hosts hundreds of influential speakers from around the world including leading scientists, renowned experts in technology, book authors, and leading academics, and makes videos of these lectures freely available.

2013 © Microsoft Corporation. All rights reserved.

NIPS Thanks Its Sponsors



amazon.com

Microsoft
Research

Google

facebook

SKYTREE.
THE MACHINE LEARNING COMPANY

TWO²⁰SIGMA

United Technologies
Research Center

YAHOO!
LABS

IBM
Research

xerox

DE Shaw & Co



DRW TRADING GROUP

criteo

PDT PARTNERS

Springer
Machine Learning Journal

Mickey Mouse icon
Disney Research

millionshort

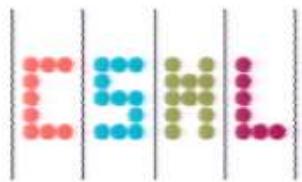
A Kernel Test for Three Variable Interactions

Dino Sejdinovic¹, Arthur Gretton¹, Wicher Bergsma²

¹Gatsby Unit, CSML, University College London

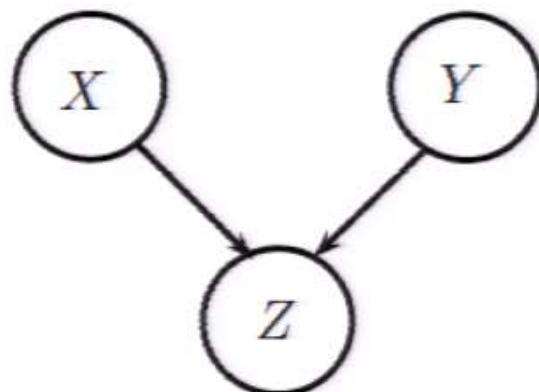
²Department of Statistics, London School of Economics

NIPS, 07 Dec 2013



Detecting a higher order interaction

- How to detect V-structures with pairwise weak (or nonexistent) dependence?



Detecting a higher order interaction

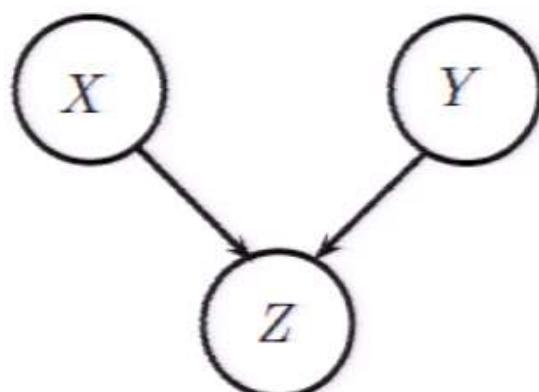
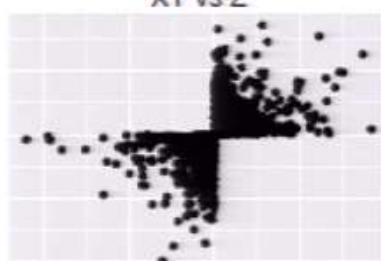
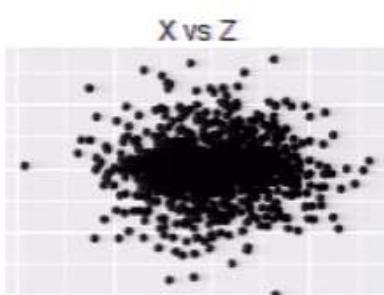
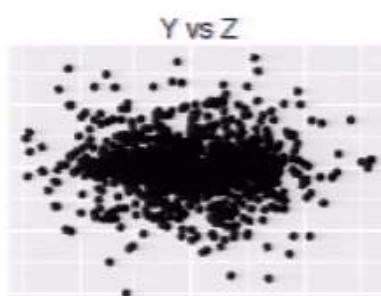
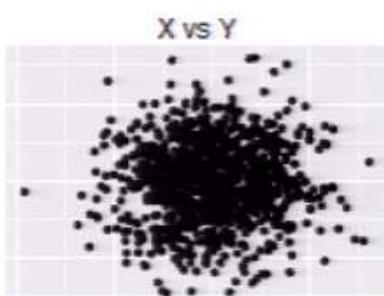
- How to detect V-structures with pairwise weak (or nonexistent) dependence?



Detecting a higher order interaction

- How to detect V-structures with pairwise weak (or nonexistent) dependence?

- $X \perp\!\!\!\perp Y, Y \perp\!\!\!\perp Z, X \perp\!\!\!\perp Z$



- $X, Y \stackrel{i.i.d.}{\sim} \mathcal{N}(0, 1)$,
- $Z | X, Y \sim \text{sign}(XY) \text{Exp}\left(\frac{1}{\sqrt{2}}\right)$

Detecting pairwise dependence

- How to detect dependence in a non-Euclidean / structured domain?

X_1 : Honourable senators, I have a question for the Leader of the Government in the Senate with regard to the support funding to farmers that has been announced. Most farmers have not received any money yet.

X_2 : No doubt there is great pressure on provincial and municipal governments in relation to the issue of child care, but the reality is that there have been no cuts to child care funding from the federal government to the provinces. In fact, we have increased federal investments for early childhood development.

...



Y_1 : Honorables sénateurs, ma question s'adresse au leader du gouvernement au Sénat et concerne l'aide financière qu'on a annoncée pour les agriculteurs. La plupart des agriculteurs n'ont encore rien reçu de cet argent.

Y_2 : Il est évident que les ordres de gouvernements provinciaux et municipaux subissent de fortes pressions en ce qui concerne les services de garde, mais le gouvernement n'a pas réduit le financement qu'il verse aux provinces pour les services de garde. Au contraire, nous avons augmenté le financement fédéral pour le développement des jeunes enfants.

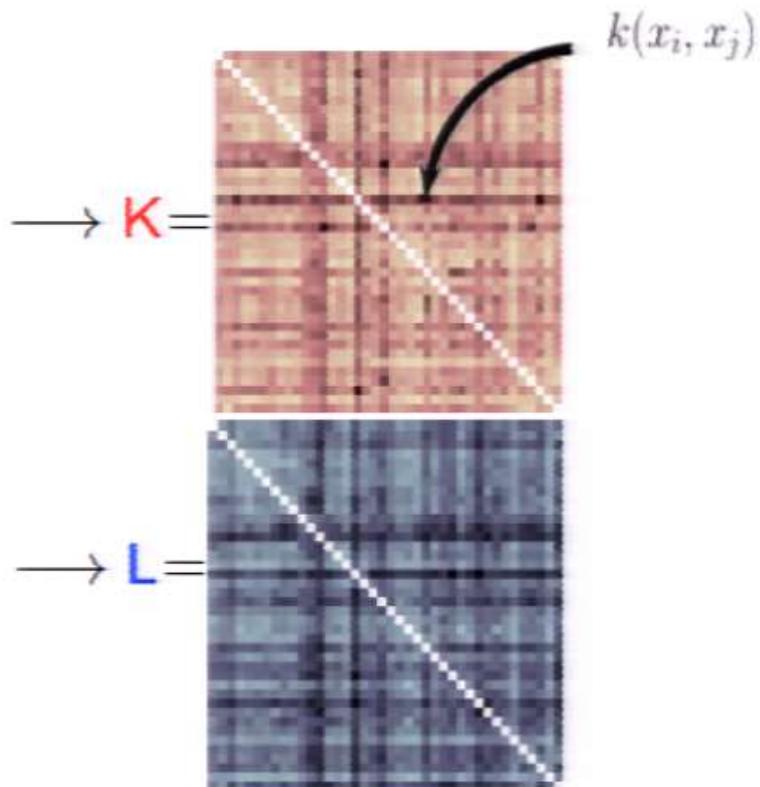
...

Are the French text extracts translations of the English ones?

Detecting pairwise dependence

.... In short, there is great pressure on provincial and municipal governments in relation to the issue of child care, but the reality is that there have been no cuts to child care funding from the federal government to the provinces. In fact, we have increased federal investments for early childhood development. ...

... il est évident que les ordres de gouvernements provinciaux et municipaux subissent de fortes pressions en ce qui concerne les services de garde, mais le gouvernement n'a pas réduit le financement qu'il verse aux provinces pour les services de garde. Au contraire, nous avons augmenté le financement fédéral pour le développement des jeunes enfants ...

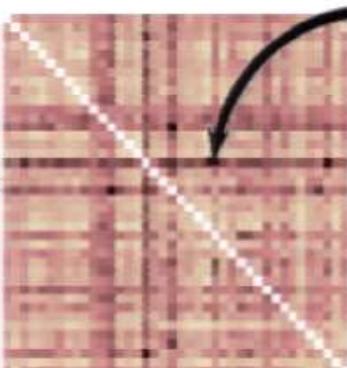


Detecting pairwise dependence

.... no doubt there is great pressure on provincial and municipal governments in relation to the issue of child care, but the reality is that there have been no cuts to child care funding from the federal government to the provinces. In fact, we have increased federal investments for early childhood development.

... il est évident que les ordres de gouvernements provinciaux et municipaux subissent de fortes pressions en ce qui concerne les services de garde, mais le gouvernement n'a pas réduit le financement qu'il verse aux provinces pour les services de garde. Au contraire, nous avons augmenté le financement fédéral pour le développement des jeunes enfants.

→ $K =$



$$k(x_i, x_j)$$

→ $L =$



- $\langle H\mathbf{K}H, H\mathbf{L}H \rangle =$

$$(H\mathbf{K}H \circ H\mathbf{L}H)_{++}$$

- $H = I - \frac{1}{n}\mathbf{1}\mathbf{1}^\top$
(centering matrix)

- $A_{++} = \sum_{i=1}^n \sum_{j=1}^n A_{ij}$

Kernel Embedding

- feature map: $z \mapsto k(\cdot, z) \in \mathcal{H}_k$
instead of $z \mapsto (\varphi_1(z), \dots, \varphi_s(z)) \in \mathbb{R}^s$
- $\langle k(\cdot, z), k(\cdot, w) \rangle_{\mathcal{H}_k} = k(z, w)$
inner products easily **computed**

Kernel Embedding

- feature map: $z \mapsto k(\cdot, z) \in \mathcal{H}_k$
instead of $z \mapsto (\varphi_1(z), \dots, \varphi_s(z)) \in \mathbb{R}^s$
- $\langle k(\cdot, z), k(\cdot, w) \rangle_{\mathcal{H}_k} = k(z, w)$
inner products easily **computed**
- embedding: $P \mapsto \mu_k(P) = \mathbb{E}_{Z \sim P} k(\cdot, Z) \in \mathcal{H}_k$
instead of $P \mapsto (\mathbb{E}\varphi_1(Z), \dots, \mathbb{E}\varphi_s(Z)) \in \mathbb{R}^s$
- $\langle \mu_k(P), \mu_k(Q) \rangle_{\mathcal{H}_k} = \mathbb{E}_{Z \sim P, W \sim Q} k(Z, W)$
inner products easily **estimated**

Independence test via embeddings

- Maximum Mean Discrepancy (MMD)

(Borgwardt et al, 2006; Gretton et al, 2007):

$$MMD_k(P, Q) = \|\mu_k(P) - \mu_k(Q)\|_{\mathcal{H}_k}$$

- ISPD kernels: μ_k injective on all signed measures and MMD_k metric (Sriperumbudur, 2010)

- Gaussian, Laplacian, inverse multiquadratics, Matérn etc.

Independence test via embeddings

- Maximum Mean Discrepancy (MMD)

(Borgwardt et al, 2006; Gretton et al, 2007):

$$MMD_k(P, Q) = \|\mu_k(P) - \mu_k(Q)\|_{\mathcal{H}_k}$$

- ISPD kernels: μ_k injective on all signed measures and MMD_k metric (Sriperumbudur, 2010)

- Gaussian, Laplacian, inverse multiquadratics, Matérn etc.

- Hilbert-Schmidt Independence Criterion (HSIC)

Gretton et al (2005, 2008); Smola et al (2007):

$$\|\mu_\kappa(P_{XY}) - \mu_\kappa(P_X P_Y)\|_{\mathcal{H}_\kappa}^2$$

$$\begin{aligned} k(\textcolor{red}{\boxed{1}}, \textcolor{red}{\boxed{2}}) &\quad l(\textcolor{blue}{\boxed{1}}, \textcolor{blue}{\boxed{2}}) \\ &\downarrow \\ \kappa(\textcolor{red}{\boxed{1}} \textcolor{blue}{\boxed{2}}, \textcolor{red}{\boxed{2}} \textcolor{blue}{\boxed{2}}) &= \\ k(\textcolor{red}{\boxed{1}}, \textcolor{red}{\boxed{2}}) \times l(\textcolor{blue}{\boxed{1}}, \textcolor{blue}{\boxed{2}}) \end{aligned}$$

Independence test via embeddings

- Maximum Mean Discrepancy (MMD)

(Borgwardt et al, 2006; Gretton et al, 2007):

$$MMD_k(P, Q) = \|\mu_k(P) - \mu_k(Q)\|_{\mathcal{H}_k}$$

- ISPD kernels: μ_k injective on all signed measures and MMD_k metric (Sriperumbudur, 2010)

- Gaussian, Laplacian, inverse multiquadratics, Matérn etc.

- Hilbert-Schmidt Independence Criterion (HSIC)

Gretton et al (2005, 2008); Smola et al (2007):

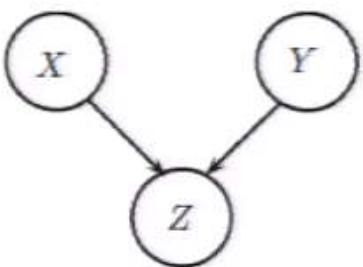
$$\|\mu_\kappa(P_{XY}) - \mu_\kappa(P_X P_Y)\|_{\mathcal{H}_\kappa}^2$$

- Empirical HSIC = $\frac{1}{n^2} \left[(HKH \circ H\mathbf{L}H)_{++} \right]$

Powerful independence tests that generalize dCov
of Szekely et al (2007); DS et al (2013)

$$\begin{aligned} k(\square_1, \square_2) &\quad l(\square_1, \square_1) \\ &\quad \downarrow \\ \kappa(\square_1 \square_2, \square_1 \square_2) &= \\ k(\square_1, \square_2) \times l(\square_1, \square_2) \end{aligned}$$

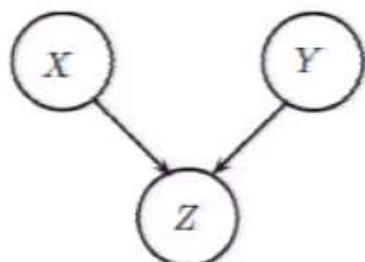
V-structure Discovery



Assume $X \perp\!\!\!\perp Y$ has been established. V-structure can then be detected by:

- CI test: $H_0 : X \perp\!\!\!\perp Y | Z$ (Zhang et al 2011) or

V-structure Discovery

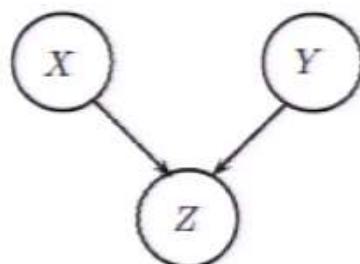
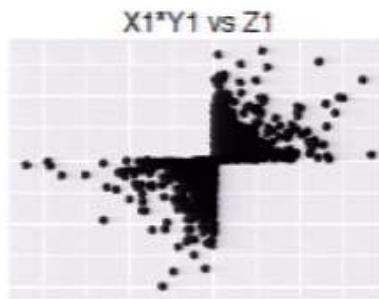
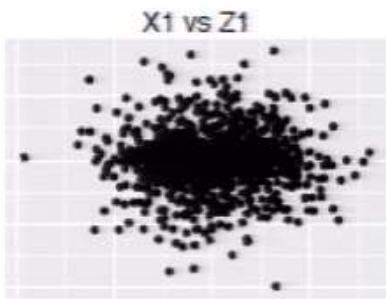
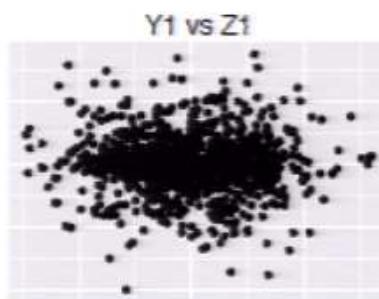
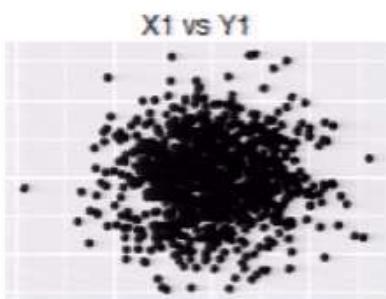


Assume $X \perp\!\!\!\perp Y$ has been established. V-structure can then be detected by:

- CI test: $H_0 : X \perp\!\!\!\perp Y | Z$ (Zhang et al 2011) or
- Factorisation test: $H_0 : (X, Y) \perp\!\!\!\perp Z \vee (X, Z) \perp\!\!\!\perp Y \vee (Y, Z) \perp\!\!\!\perp X$ (multiple standard two-variable tests)
 - compute p -values for each of the marginal tests for $(Y, Z) \perp\!\!\!\perp X$, $(X, Z) \perp\!\!\!\perp Y$, or $(X, Y) \perp\!\!\!\perp Z$
 - apply Holm-Bonferroni (**HB**) sequentially rejective correction (Holm 1979)

V-structure Discovery (2)

- How to detect V-structures with pairwise weak (or nonexistent) dependence?
- $X \perp\!\!\!\perp Y, Y \perp\!\!\!\perp Z, X \perp\!\!\!\perp Z$



- $X_1, Y_1 \stackrel{i.i.d.}{\sim} \mathcal{N}(0, 1)$,
- $Z_1 | X_1, Y_1 \sim \text{sign}(X_1 Y_1) \text{Exp}\left(\frac{1}{\sqrt{2}}\right)$

V-structure Discovery (3)

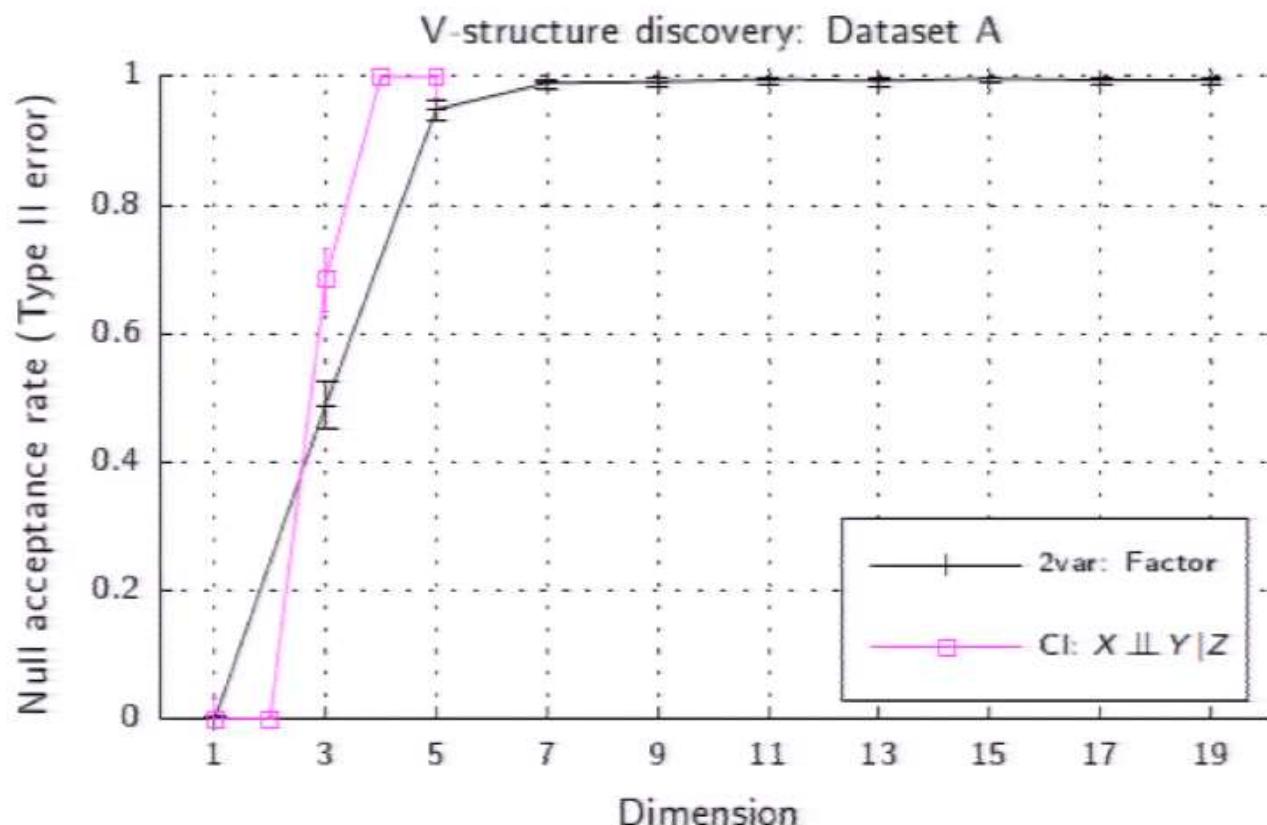


Figure: CI test for $X \perp\!\!\! \perp Y | Z$ from Zhang et al (2011), and a factorisation test with a HB correction, $n = 500$

Lancaster Interaction Measure

Definition (Bahadur (1961); Lancaster (1969))

Interaction measure of $(X_1, \dots, X_D) \sim P$ is a signed measure ΔP that **vanishes** whenever P can be factorised in a non-trivial way as a product of its (possibly multivariate) marginal distributions.

Lancaster Interaction Measure

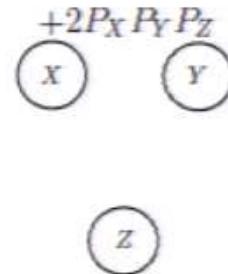
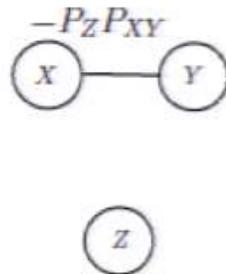
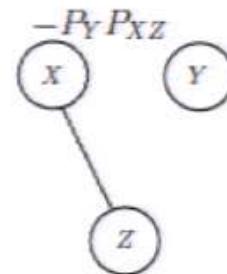
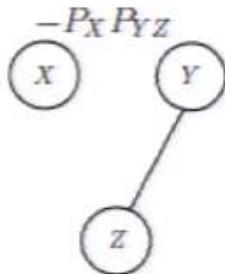
Definition (Bahadur (1961); Lancaster (1969))

Interaction measure of $(X_1, \dots, X_D) \sim P$ is a signed measure ΔP that vanishes whenever P can be factorised in a non-trivial way as a product of its (possibly multivariate) marginal distributions.

- $D = 2 : \Delta_L P = P_{XY} - P_X P_Y$
- $D = 3 : \Delta_L P = P_{XYZ} - P_X P_{YZ} - P_Y P_{XZ} - P_Z P_{XY} + 2P_X P_Y P_Z$

$$\Delta_L P =$$

$$P_{XYZ}$$



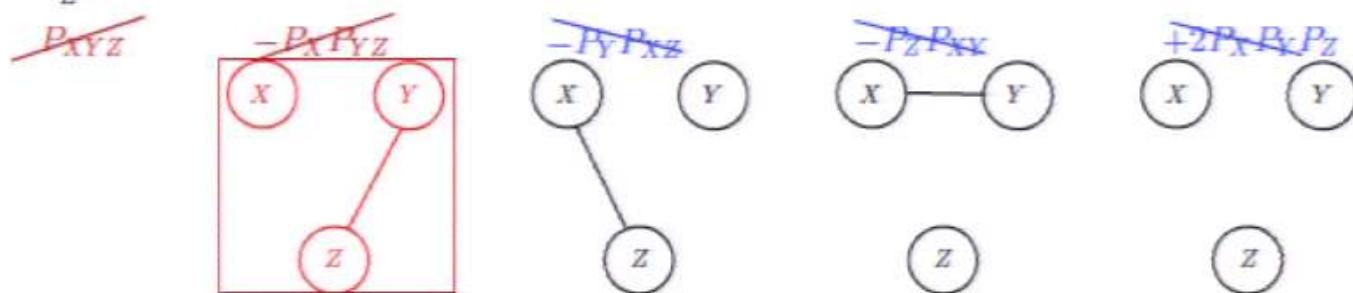
Lancaster Interaction Measure

Definition (Bahadur (1961); Lancaster (1969))

Interaction measure of $(X_1, \dots, X_D) \sim P$ is a signed measure ΔP that vanishes whenever P can be factorised in a non-trivial way as a product of its (possibly multivariate) marginal distributions.

- $D = 2 : \Delta_L P = P_{XY} - P_X P_Y$
- $D = 3 : \Delta_L P = P_{XYZ} - P_X P_{YZ} - P_Y P_{XZ} - P_Z P_{XY} + 2P_X P_Y P_Z$

$$\Delta_L P = 0$$



A Test using Lancaster Measure

- Construct a test by estimating $\|\mu_\kappa(\Delta_L P)\|_{\mathcal{H}_\kappa}^2$, where $\kappa = \mathbf{k} \otimes \mathbf{l} \otimes \mathbf{m}$:

$$\begin{aligned}\|\mu_\kappa(P_{XYZ} - P_{XY}P_Z - \dots)\|_{\mathcal{H}_\kappa}^2 &= \\ \langle \mu_\kappa P_{XYZ}, \mu_\kappa P_{XYZ} \rangle_{\mathcal{H}_\kappa} - 2 \langle \mu_\kappa P_{XYZ}, \mu_\kappa P_{XY}P_Z \rangle_{\mathcal{H}_\kappa} - \dots\end{aligned}$$

Inner Product Estimators

$\nu \backslash \nu'$	P_{XYZ}	$P_{XY}P_Z$	$P_{XZ}P_Y$	$P_{YZ}P_X$	$P_XP_YP_Z$
P_{XYZ}	$(K \circ L \circ M)_{++}$	$((K \circ L) M)_{++}$	$((K \circ M) L)_{++}$	$((M \circ L) K)_{++}$	$tr(K_+ \circ L_+ \circ M_+)$
$P_{XY}P_Z$		$(K \circ L)_{++} M_{++}$	$(MKL)_{++}$	$(KLM)_{++}$	$(KL)_{++} M_{++}$
$P_{XZ}P_Y$			$(K \circ M)_{++} L_{++}$	$(KML)_{++}$	$(KM)_{++} L_{++}$
$P_{YZ}P_X$				$(L \circ M)_{++} K_{++}$	$(LM)_{++} K_{++}$
$P_XP_YP_Z$					$K_{++} L_{++} M_{++}$

Table: V -statistic estimators of $\langle \mu_\kappa \nu, \mu_\kappa \nu' \rangle_{\mathcal{H}_\kappa}$

Inner Product Estimators

$\nu \backslash \nu'$	P_{XYZ}	$P_{XY}P_Z$	$P_{XZ}P_Y$	$P_{YZ}P_X$	$P_XP_YP_Z$
P_{XYZ}	$(K \circ L \circ M)_{++}$	$((K \circ L) M)_{++}$	$((K \circ M) L)_{++}$	$((M \circ L) K)_{++}$	$tr(K_+ \circ L_+ \circ M_+)$
$P_{XY}P_Z$		$(K \circ L)_{++} M_{++}$	$(M K L)_{++}$	$(K L M)_{++}$	$(K L)_{++} M_{++}$
$P_{XZ}P_Y$			$(K \circ M)_{++} L_{++}$	$(K M L)_{++}$	$(K M)_{++} L_{++}$
$P_{YZ}P_X$				$(L \circ M)_{++} K_{++}$	$(L M)_{++} K_{++}$
$P_XP_YP_Z$					$K_{++} L_{++} M_{++}$

Table: V -statistic estimators of $\langle \mu_\kappa \nu, \mu_\kappa \nu' \rangle_{\mathcal{H}_\kappa}$

Proposition (Lancaster interaction statistic)

$$\| \mu_\kappa (\Delta_L P) \|_{\mathcal{H}_\kappa}^2 = \frac{1}{n^2} \boxed{(H \textcolor{red}{K} H \circ H \textcolor{blue}{L} H \circ H \textcolor{magenta}{M} H)_{++}}.$$

Empirical joint central moment in the feature space

Example A: factorisation tests

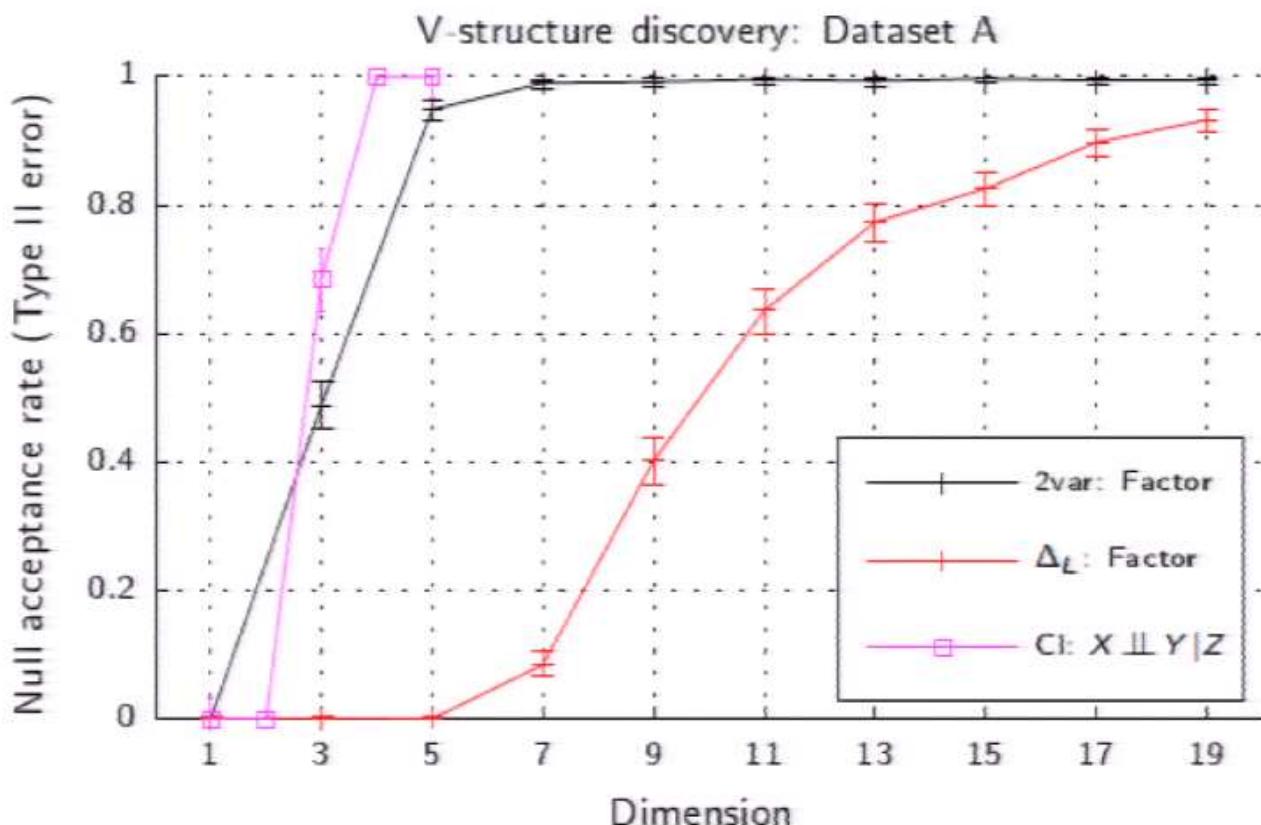


Figure: Factorisation hypothesis: Lancaster statistic vs. a two-variable based test (both with HB correction); Test for $X \perp\!\!\!\perp Y|Z$ from Zhang et al (2011), $n = 500$

Example B: Joint dependence can be easier to detect

- $X_1, Y_1 \stackrel{i.i.d.}{\sim} \mathcal{N}(0, 1)$
- $Z_1 = \begin{cases} X_1^2 + \epsilon, & w.p. 1/3, \\ Y_1^2 + \epsilon, & w.p. 1/3, \\ X_1 Y_1 + \epsilon, & w.p. 1/3, \end{cases}$ where $\epsilon \sim \mathcal{N}(0, 0.1^2)$.
- $X_{2:p}, Y_{2:p}, Z_{2:p} \stackrel{i.i.d.}{\sim} \mathcal{N}(0, I_{p-1})$
- dependence of Z on pair (X, Y) is stronger than on X and Y individually

Example B: factorisation tests

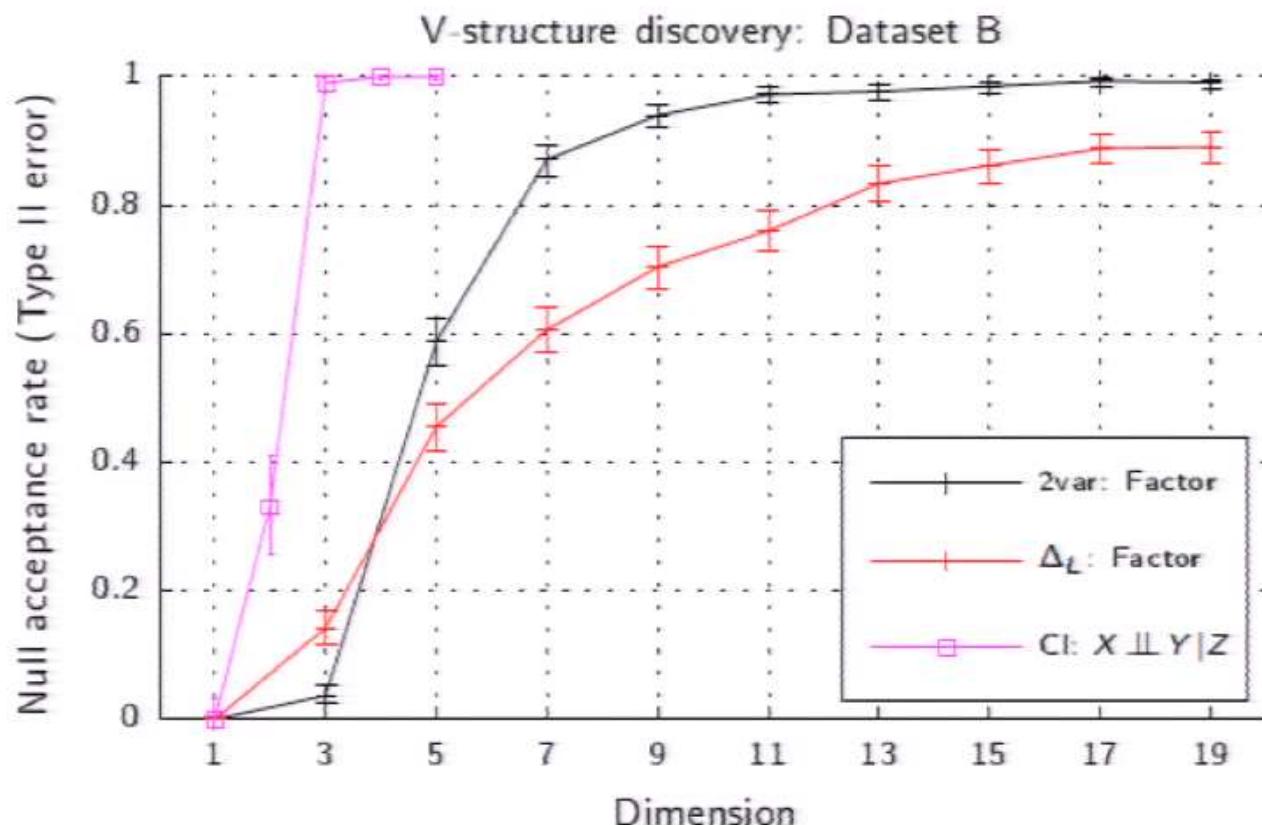


Figure: Factorisation hypothesis: Lancaster statistic vs. a two-variable based test (both with HB correction); Test for $X \perp\!\!\! \perp Y|Z$ from Zhang et al (2011), $n = 500$

Interaction for $D \geq 4$

- Interaction measure valid for all D
(Streitberg, 1990):

$$\Delta_S P = \sum_{\pi} (-1)^{|\pi|-1} (|\pi| - 1)! J_{\pi} P$$

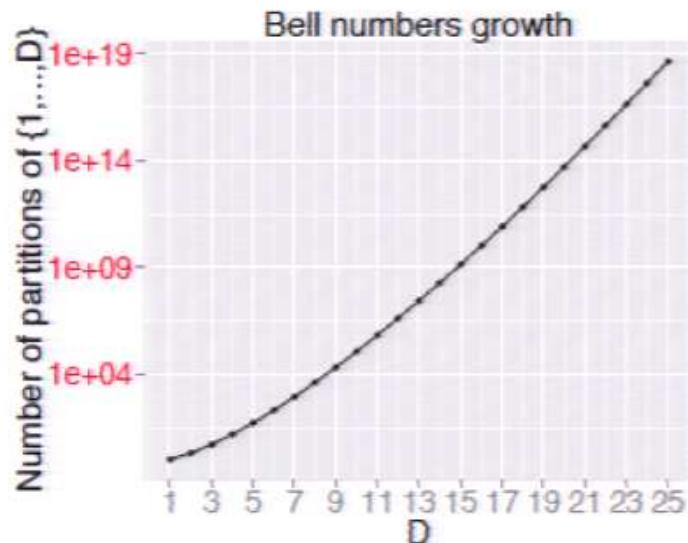
- For a partition π , J_{π} associates to the joint the corresponding factorisation,
e.g., $J_{13|2|4} P = P_{X_1 X_3} P_{X_2} P_{X_4}$.

Interaction for $D \geq 4$

- Interaction measure valid for all D
(Streitberg, 1990):

$$\Delta_S P = \sum_{\pi} (-1)^{|\pi|-1} (|\pi| - 1)! J_{\pi} P$$

- For a partition π , J_{π} associates to the joint the corresponding factorisation, e.g., $J_{13|2|4} P = P_{X_1 X_3} P_{X_2} P_{X_4}$.



Interaction for $D \geq 4$

- Interaction measure valid for all D
(Streitberg, 1990):

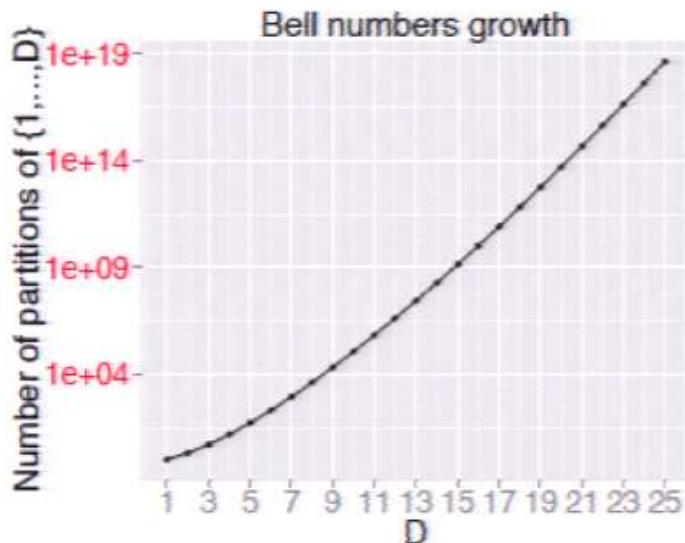
$$\Delta_S P = \sum_{\pi} (-1)^{|\pi|-1} (|\pi| - 1)! J_{\pi} P$$

- For a partition π , J_{π} associates to the joint the corresponding factorisation, e.g., $J_{13|2|4} P = P_{X_1 X_3} P_{X_2} P_{X_4}$.

joint central moments (Lancaster interaction)

vs.

joint cumulants (Streitberg interaction)



Summary

- A nonparametric test for three-variable interaction and for total independence, using embeddings of signed measures into RKHSs

Summary

- A nonparametric test for three-variable interaction and for total independence, using embeddings of signed measures into RKHSs
- Test statistics are simple and easy to compute - corresponding permutation tests significantly outperform standard two-variable-based tests on V-structures with weak pairwise interactions

Summary

- A nonparametric test for three-variable interaction and for total independence, using embeddings of signed measures into RKHSs
- Test statistics are simple and easy to compute - corresponding permutation tests significantly outperform standard two-variable-based tests on V-structures with weak pairwise interactions
- All forms of Lancaster three-variable interaction can be detected for a large family of reproducing kernels (**ISPD**)

Summary

- A nonparametric test for three-variable interaction and for total independence, using embeddings of signed measures into RKHSs
- Test statistics are simple and easy to compute - corresponding permutation tests significantly outperform standard two-variable-based tests on V-structures with weak pairwise interactions
- All forms of Lancaster three-variable interaction can be detected for a large family of reproducing kernels (**ISPD**)

Thank You!
Poster **S6**

NIPS Thanks Its Sponsors



amazon.com

Microsoft
Research

Google

facebook

SKYTREE.
THE MACHINE LEARNING COMPANY

TWO²⁰SIGMA

United Technologies
Research Center

YAHOO!
LABS

IBM
Research

xerox

DE Shaw & Co

millionshort



DRW TRADING GROUP

criteo

PDT PARTNERS

Springer
Machine Learning Journal

Mickey Mouse icon
Disney Research

More Effective Distributed ML via a Stale Synchronous Parallel Parameter Server

Q. Ho, J. Cipar, H. Cui, J.K. Kim, S. Lee,

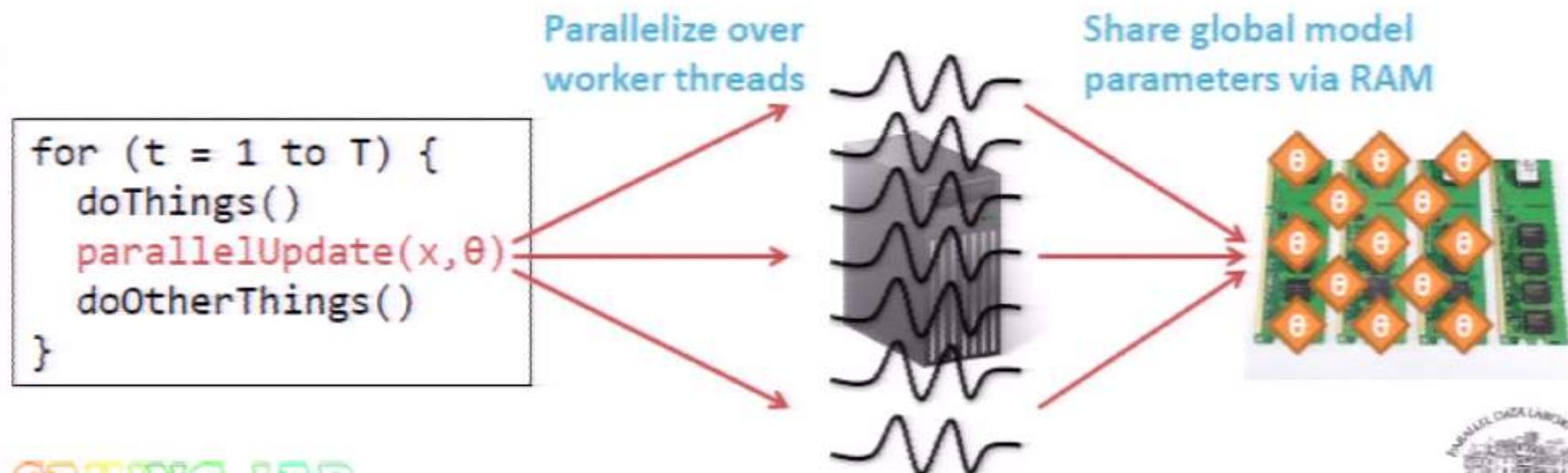
**P.B. Gibbons, G.A. Gibson, G.R. Ganger, E.P. Xing*

Carnegie Mellon University

*Intel Labs

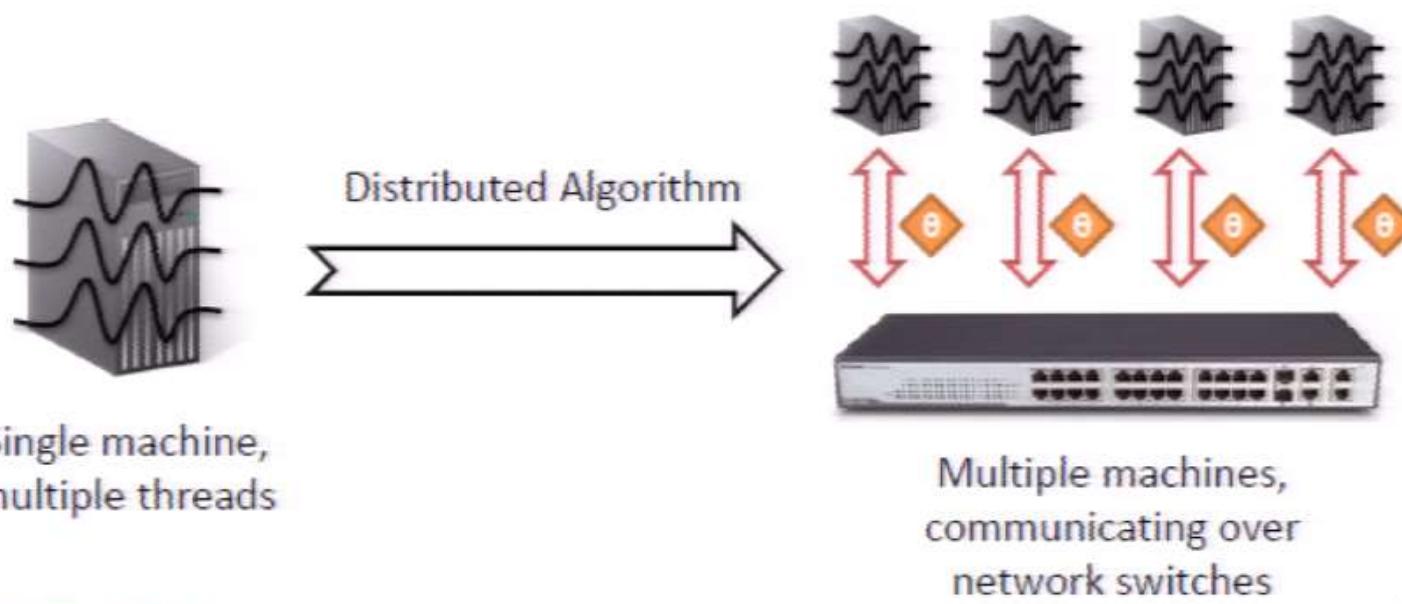
Distributed ML: one machine to many

- Setting: have **iterative, parallel ML algorithm**
 - E.g. optimization, MCMC algorithms
 - For topic models, regression, matrix factorization, SVMs, DNNs, etc.
- Critical updates executed on one machine, in parallel
 - Worker threads **share global model parameters θ via RAM**



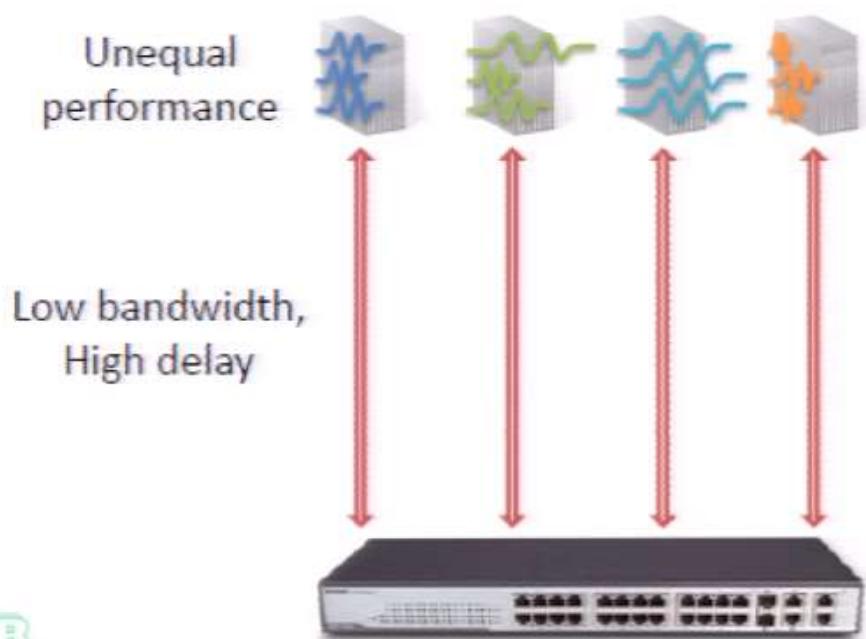
Distributed ML: one machine to many

- Want: scale up by distributing ML algorithm
 - Must now share parameters over a network
- Seems like a simple task...
 - Many distributed tools available, so just pick one and go?



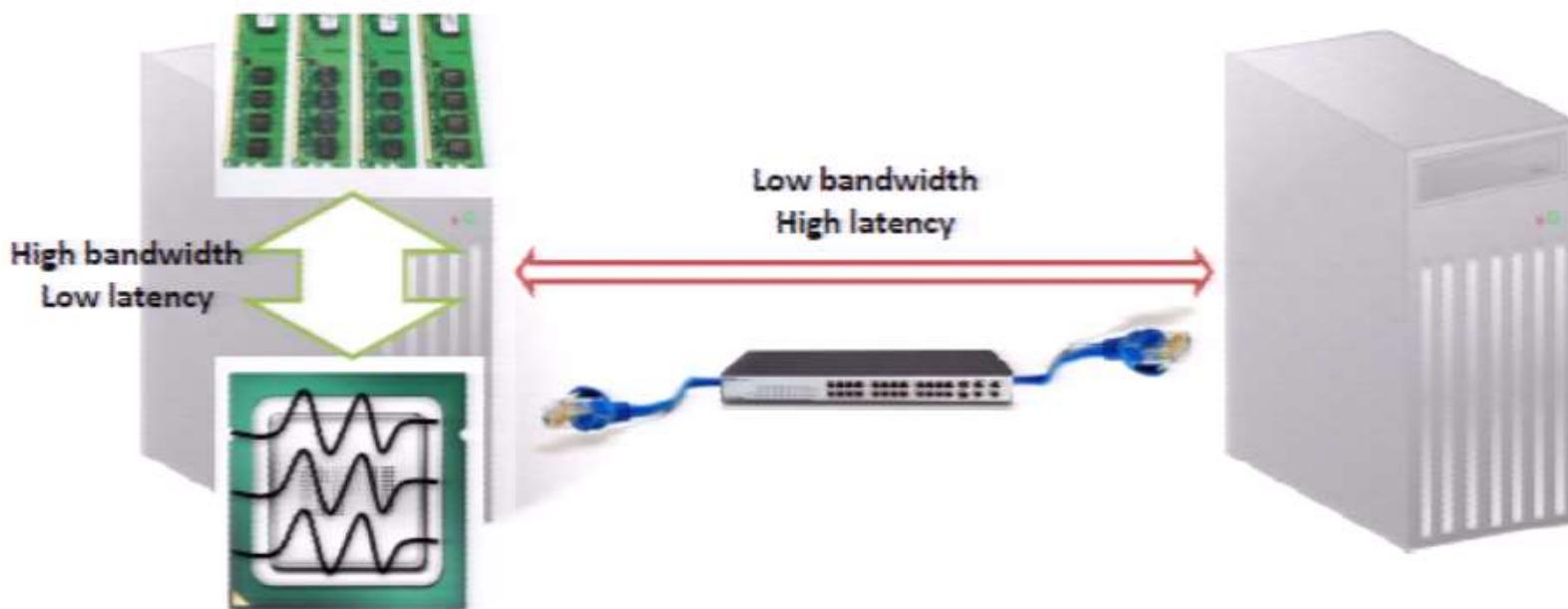
Distributed ML Challenges

- Not quite that easy...
- **Two distributed challenges:**
 - Networks are slow
 - “Identical” machines rarely perform equally



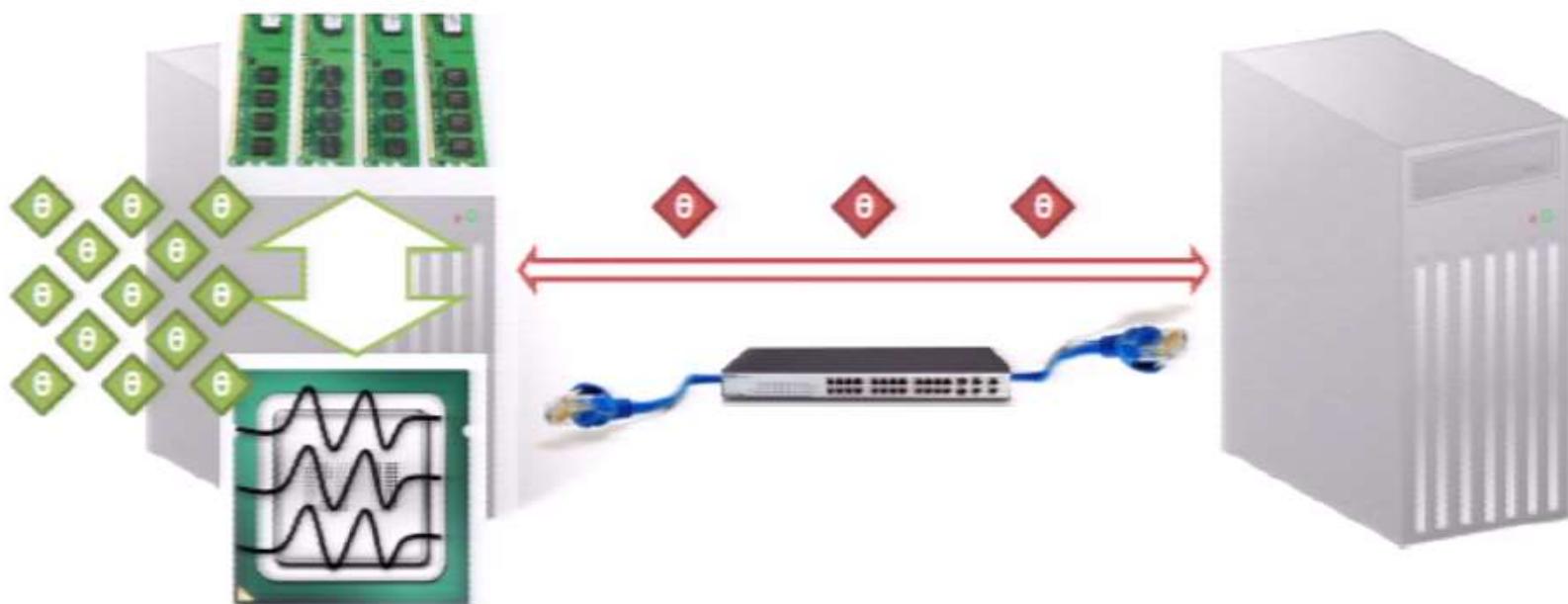
Networks are (relatively) slow

- Low network bandwidth:
 - 0.1-1GB/s (inter-machine) vs $\geq 20\text{GB/s}$ (CPU-RAM)
 - Fewer parameters transmitted per second
- High network latency (messaging time):
 - 10,000-100,000 ns (inter-machine) vs 100 ns (CPU-RAM)
 - Wait much longer to receive parameters

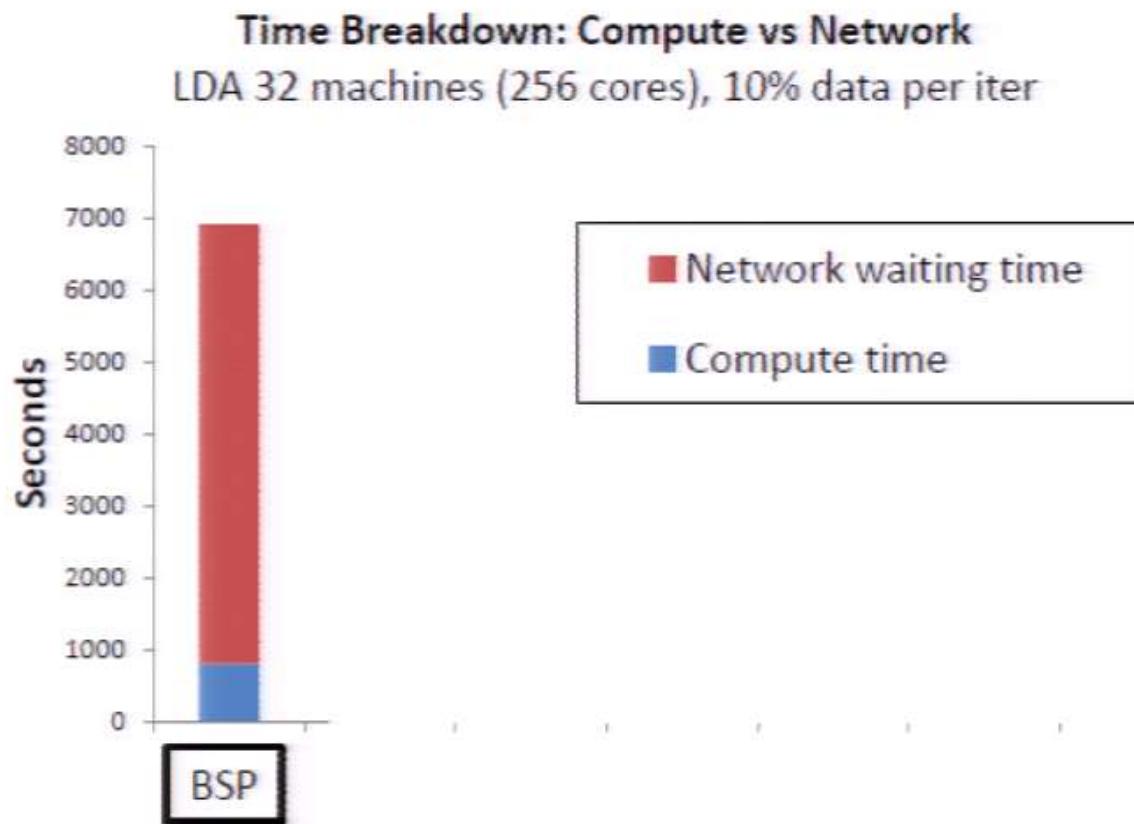


Networks are (relatively) slow

- Parallel ML requires frequent synchronization
 - Exchange 10-1000K scalars per second, per thread
 - Parameters not shared quickly enough → communication bottleneck
- Significant bottleneck over a network!



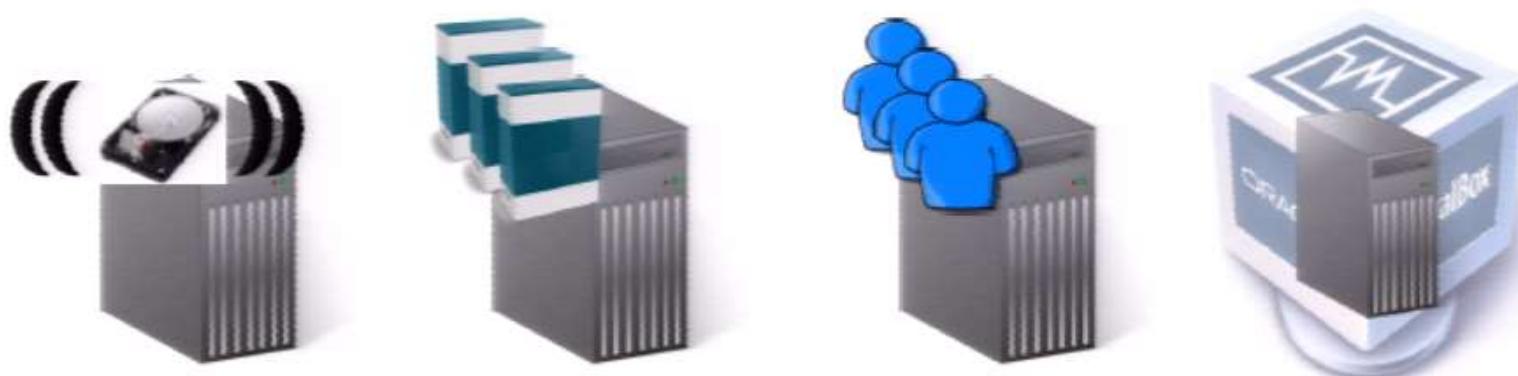
Networks are (relatively) slow



For a “clean” setting with full control over machines and full network capacity
Real clusters with many users have even worse network:compute ratios!

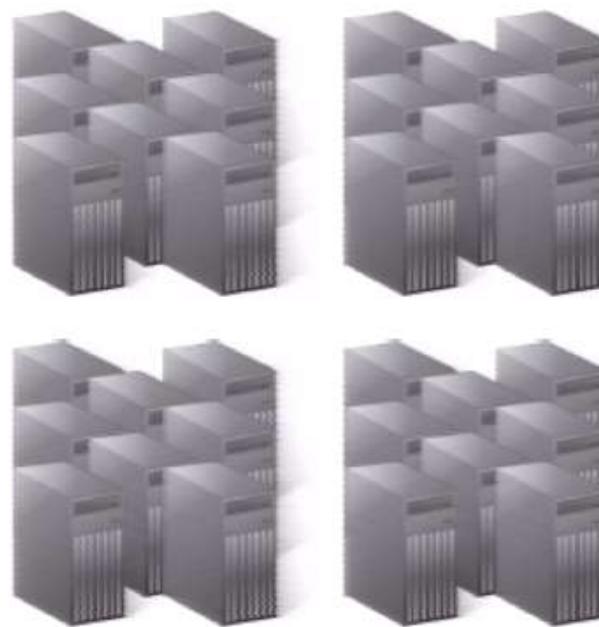
Machines don't perform equally

- Even when configured identically
- **Variety of reasons:**
 - Vibrating hard drive
 - Background programs; part of a distributed filesystem
 - Other users
 - Machine is a VM/cloud service
- **Occasional, random slowdowns in different machines**



Consequence: Scaling up ML is hard!

- **Going from 1 to N machines:**
 - Naïve implementations rarely yield N-fold speedup
 - Slower convergence due to machine slowdowns, network bottlenecks
 - If not careful, even worse than a single machine!
 - Algorithm diverges due to errors from slowdowns!



Existing general-purpose scalable ML

Theory-oriented

- Focus on algorithm correctness/convergence
- Examples:
 - Cyclic fixed-delay schemes (Langford et al., Agarwal & Duchi)
 - Single-machine asynchronous (Niu et al.)
 - Naively-parallel SGD (Zinkevich et al.)
 - Partitioned SGD (Gemulla et al.)
- May oversimplify systems issues
 - e.g. need machines to perform consistently
 - e.g. need lots of synchronization
 - e.g. or even try not to communicate at all

Systems-oriented

- Focus on high iteration throughput
- Examples:
 - MapReduce: Hadoop and Mahout
 - Spark
 - Graph-based: GraphLab, Pregel
- May oversimplify ML issues
 - e.g. assume algorithms "just work" in distributed setting, without proof
 - e.g. must convert programs to new programming model; nontrivial effort

Existing general-purpose scalable ML

Theory-oriented

- Focus on algorithm correctness/convergence
- Examples:
 - Cyclic fixed-delay schemes (Langford et al., Agarwal & Duchi)
 - Single-machine asynchronous (Niu et al.)
 - Naively-parallel SGD (Zinkevich et al.)
 - Partitioned SGD (Gemulla et al.)
- May oversimplify systems issues
 - e.g. need machines to perform consistently
 - e.g. need lots of synchronization
 - e.g. or even try not to communicate at all

Systems-oriented

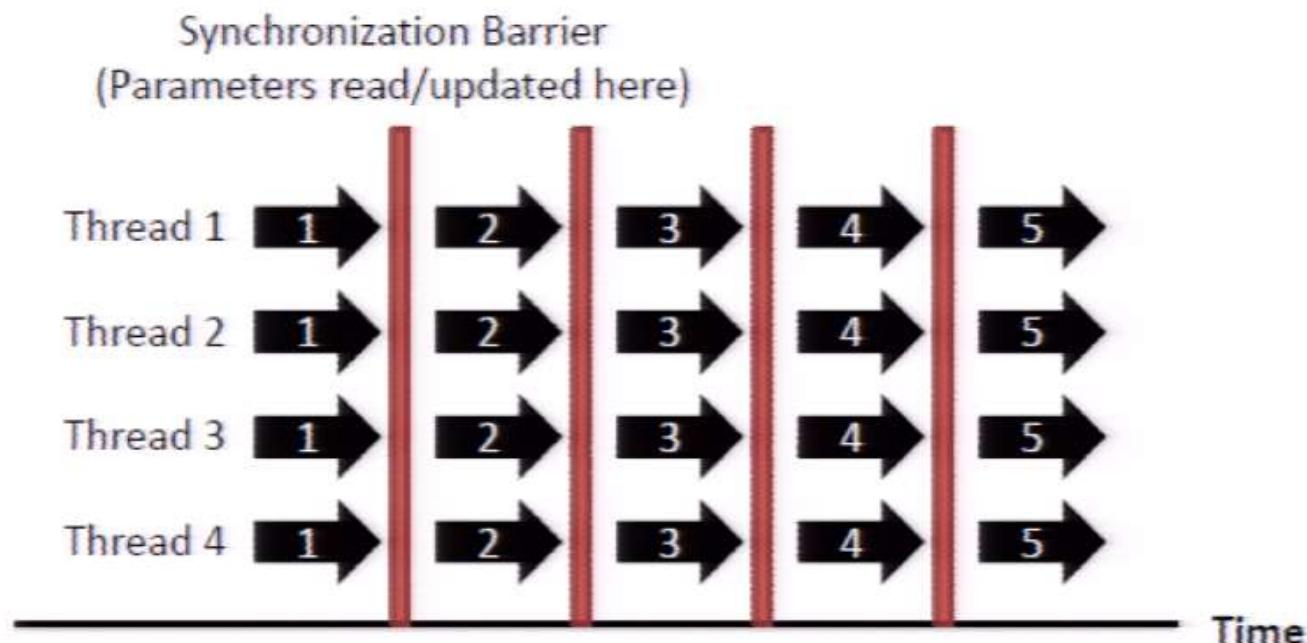
- Focus on high iteration throughput
- Examples:
 - MapReduce: Hadoop and Mahout
 - Spark
 - Graph-based: GraphLab, Pregel
- May oversimplify ML issues
 - e.g. assume algorithms “just work” in distributed setting, without proof
 - e.g. must convert programs to new programming model; nontrivial effort

Can we take both sides into account?

Middle of the road approach

- Want: ML algorithms converge quickly under **imperfect systems conditions**
 - e.g. slow network performance
 - e.g. random machine slowdowns
 - Parameters are not communicated consistently
- Existing work: mostly use one of two communication models
 - Bulk Synchronous Parallel (BSP)
 - Asynchronous (Async)
- First, understand pros and cons of BSP and Async

Bulk Synchronous Parallel



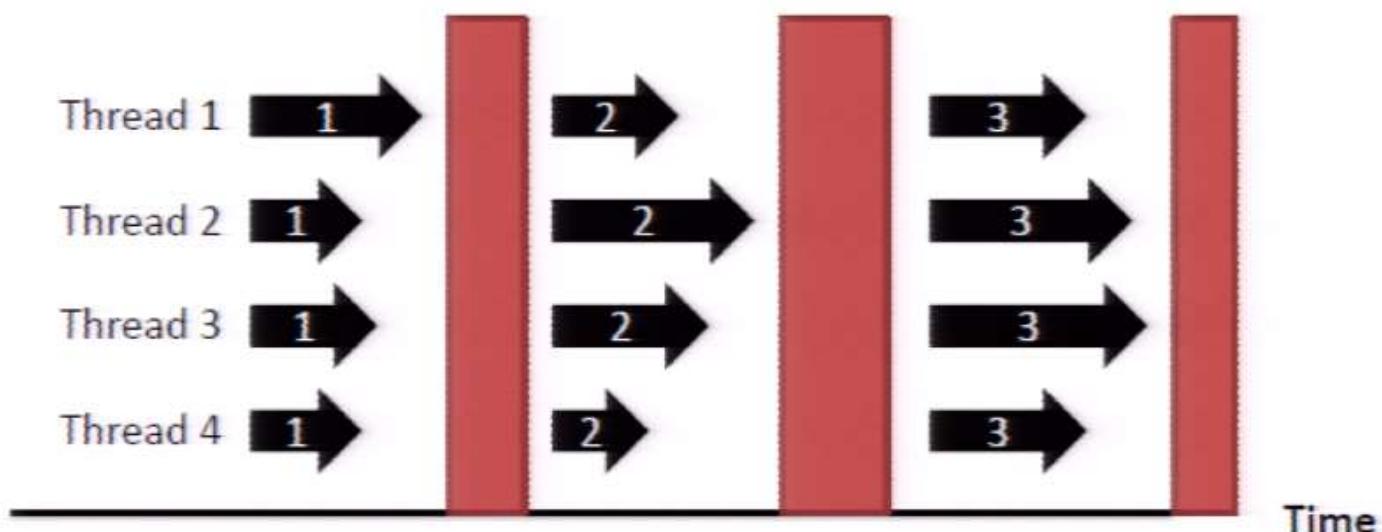
Threads synchronize (wait for each other) every iteration

Threads all on same iteration #

Parameters read/updated at synchronization barriers



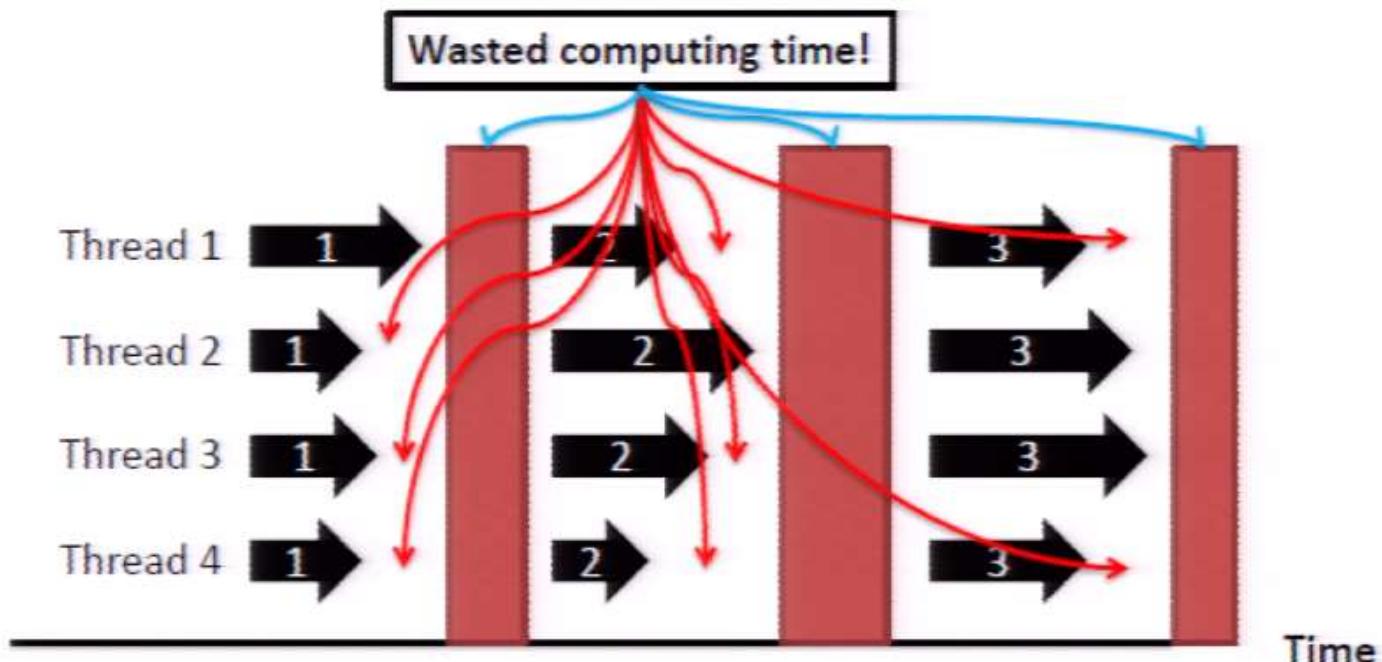
The cost of synchronicity



- (a) Machines perform unequally
 - (b) Algorithmic workload imbalanced
- So threads must wait for each other**

End-of-iteration sync gets longer with larger clusters (due to slow network)

The cost of synchronicity

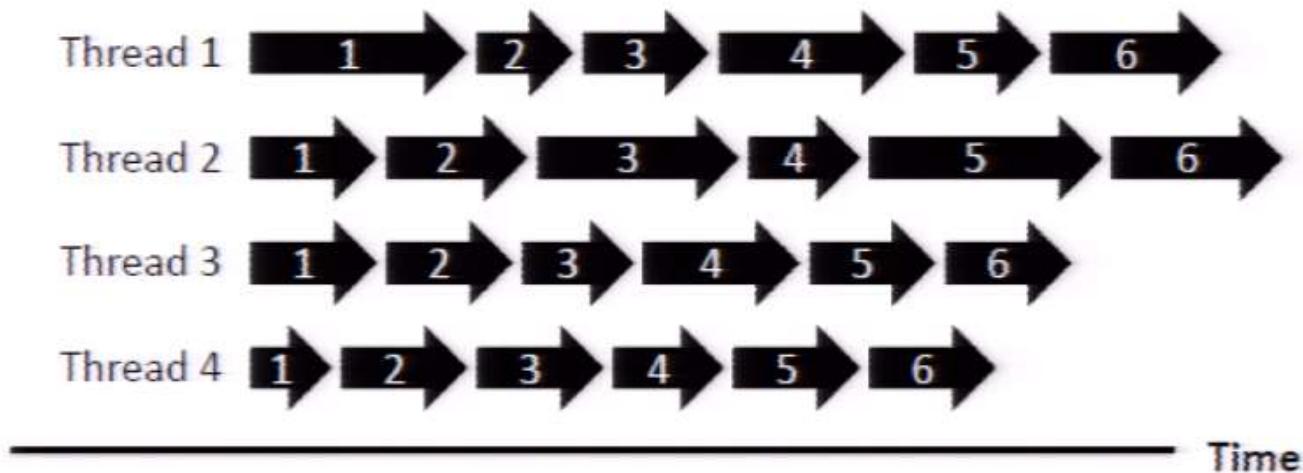


Threads must wait for each other
 End-of-iteration sync gets longer with larger clusters

Precious computing time wasted

Asynchronous

Parameters read/updated
at any time

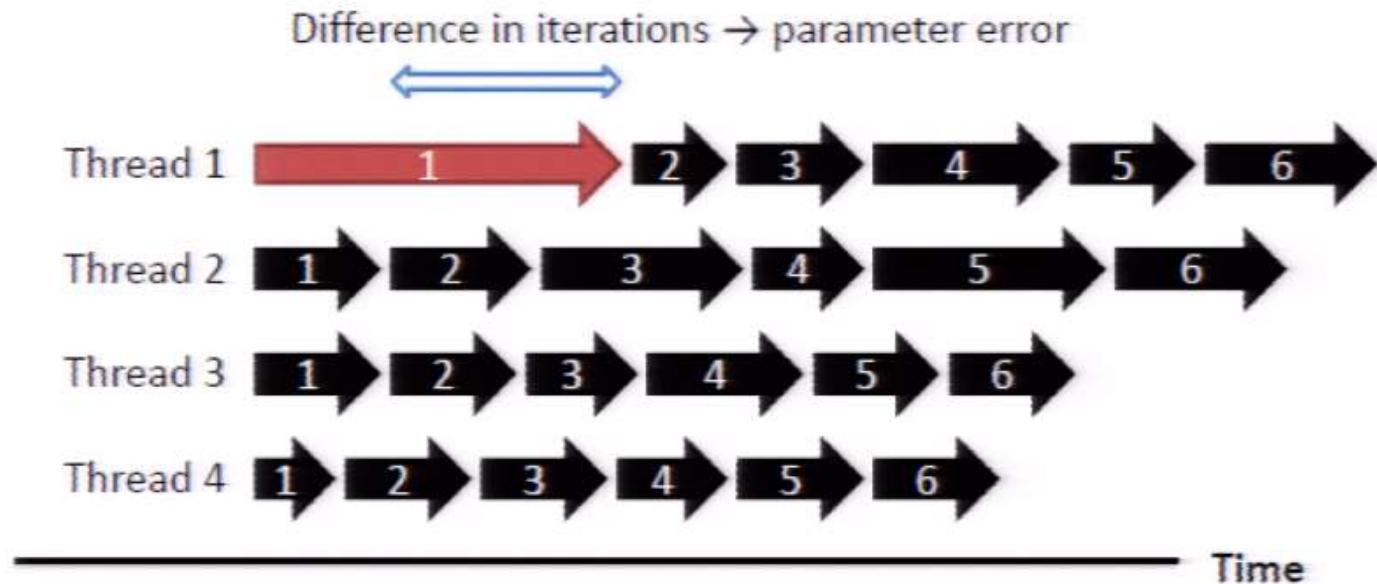


Threads proceed to next iteration without waiting

Threads not on same iteration #

Parameters read/updated any time

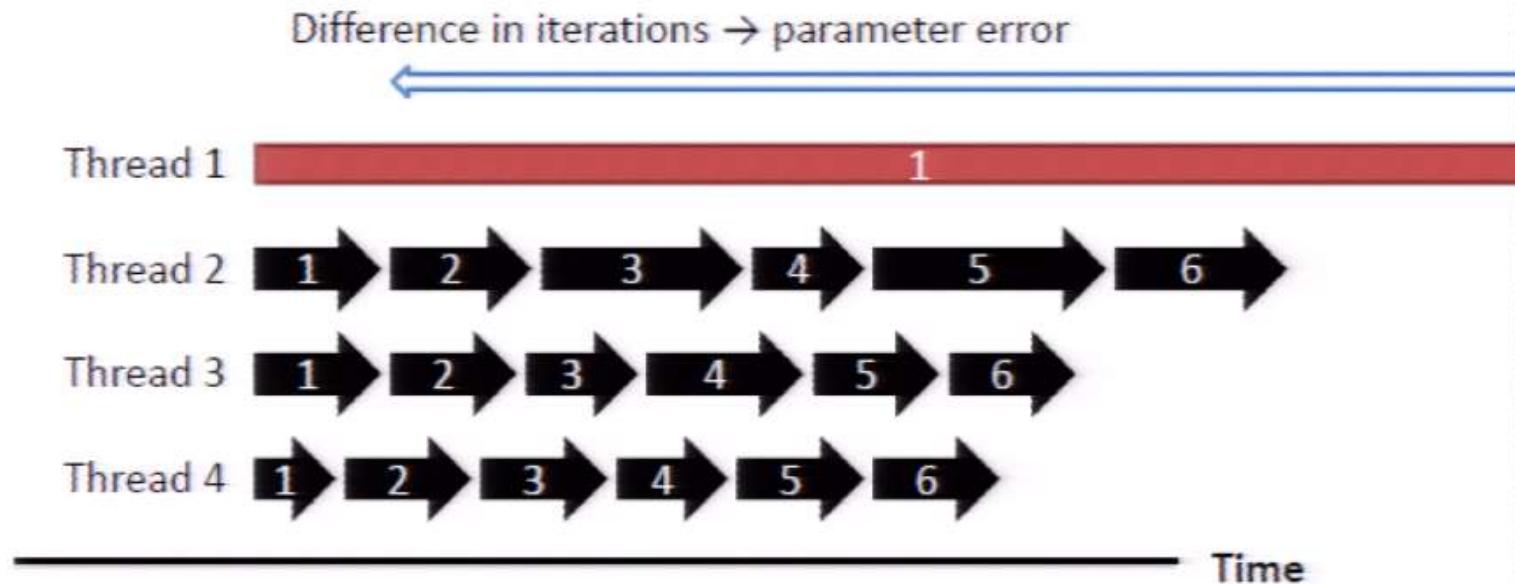
Slowdowns and Async



Machine suddenly slows down (hard drive, background process, etc.)
 Causing iteration difference between threads
 Leading to error in parameters



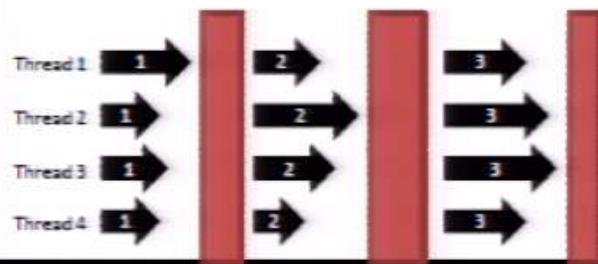
Async worst-case situation



Large clusters have arbitrarily large slowdowns!
 Machines become inaccessible for extended periods
Error becomes unbounded!

What we really want

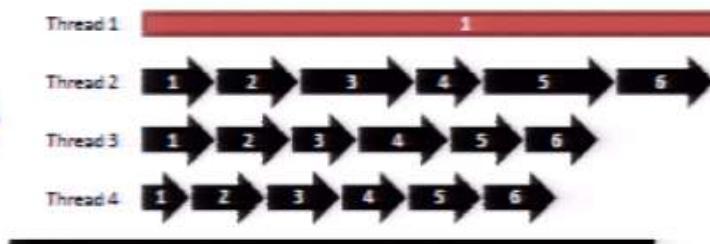
- “Partial” synchronicity
 - Spread network comms evenly (don’t sync unless needed)
 - Threads usually shouldn’t wait – but mustn’t drift too far apart!
- Straggler tolerance
 - Slow threads must somehow catch up
- Is there a middle ground between BSP and Async?



BSP



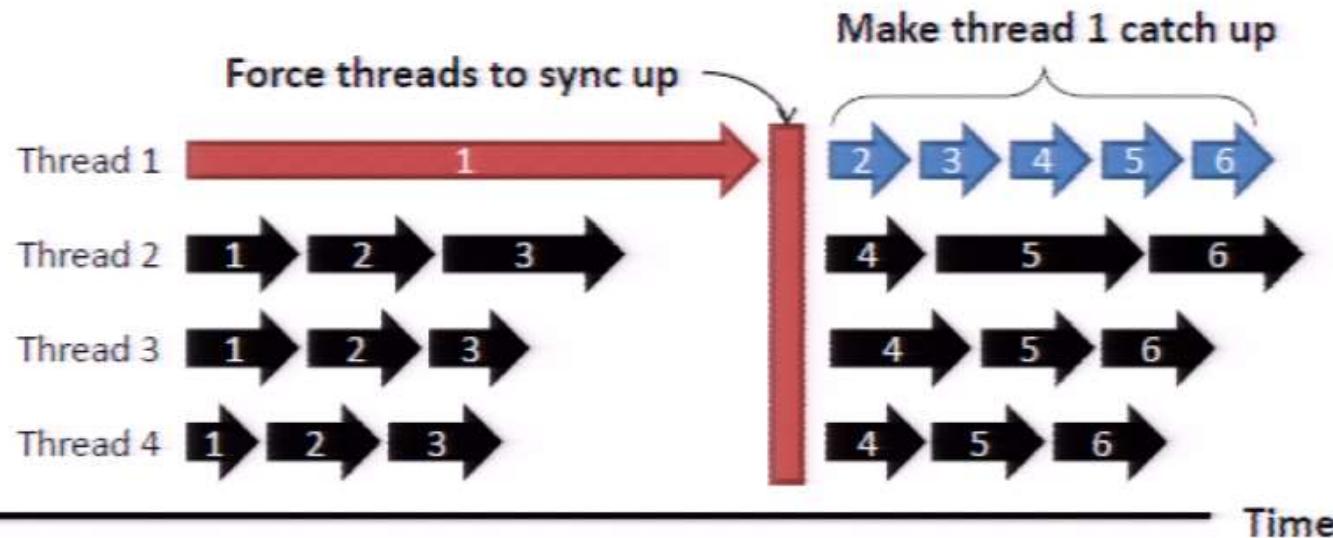
???



Async

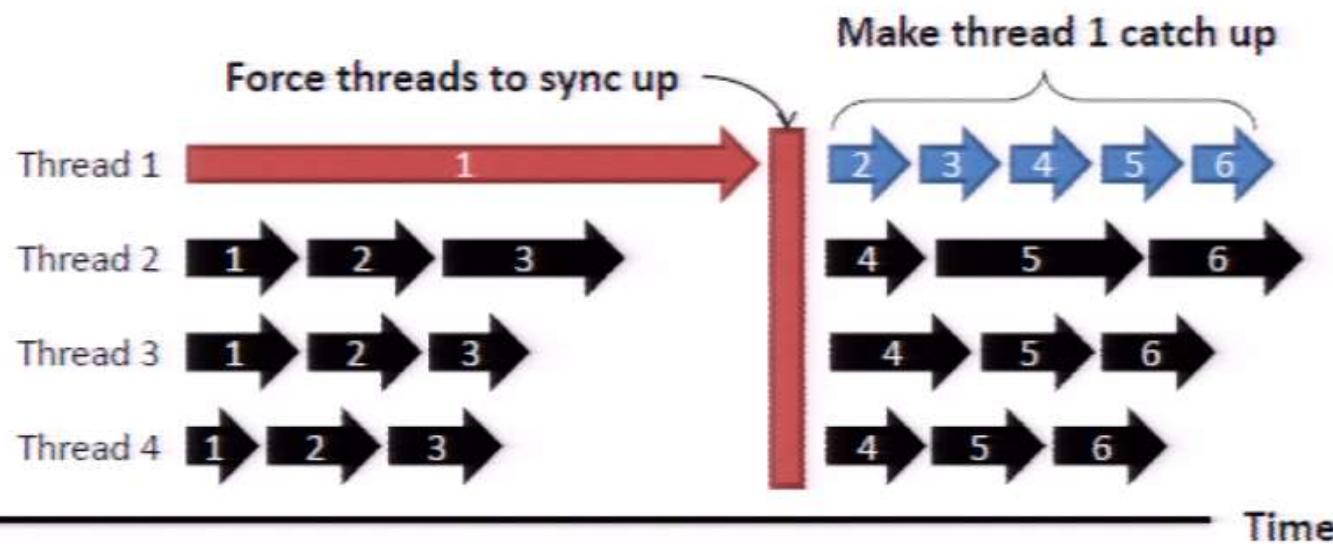
That middle ground

- “Partial” synchronicity
 - Spread network comms evenly (don’t sync unless needed)
 - Threads usually shouldn’t wait – but mustn’t drift too far apart!
- Straggler tolerance
 - Slow threads must somehow catch up

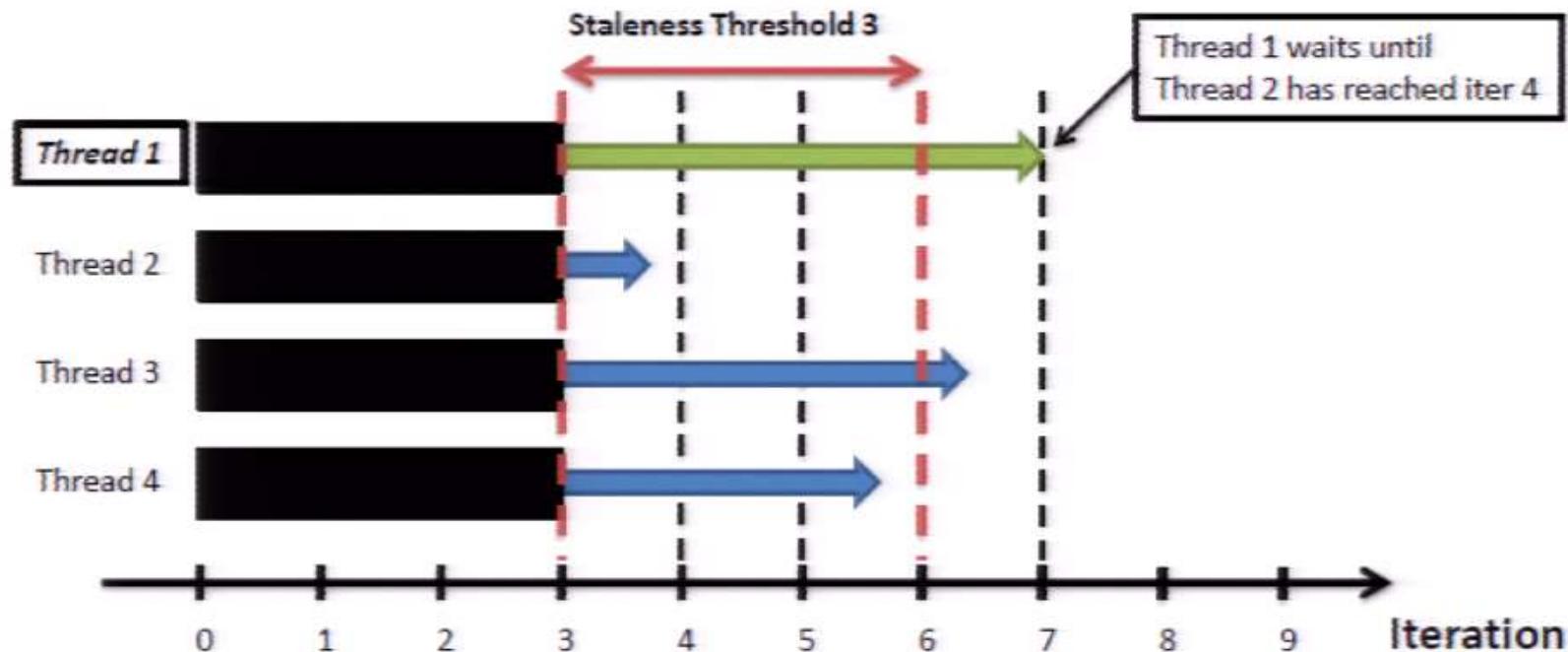


That middle ground

How do we realize this?



Stale Synchronous Parallel



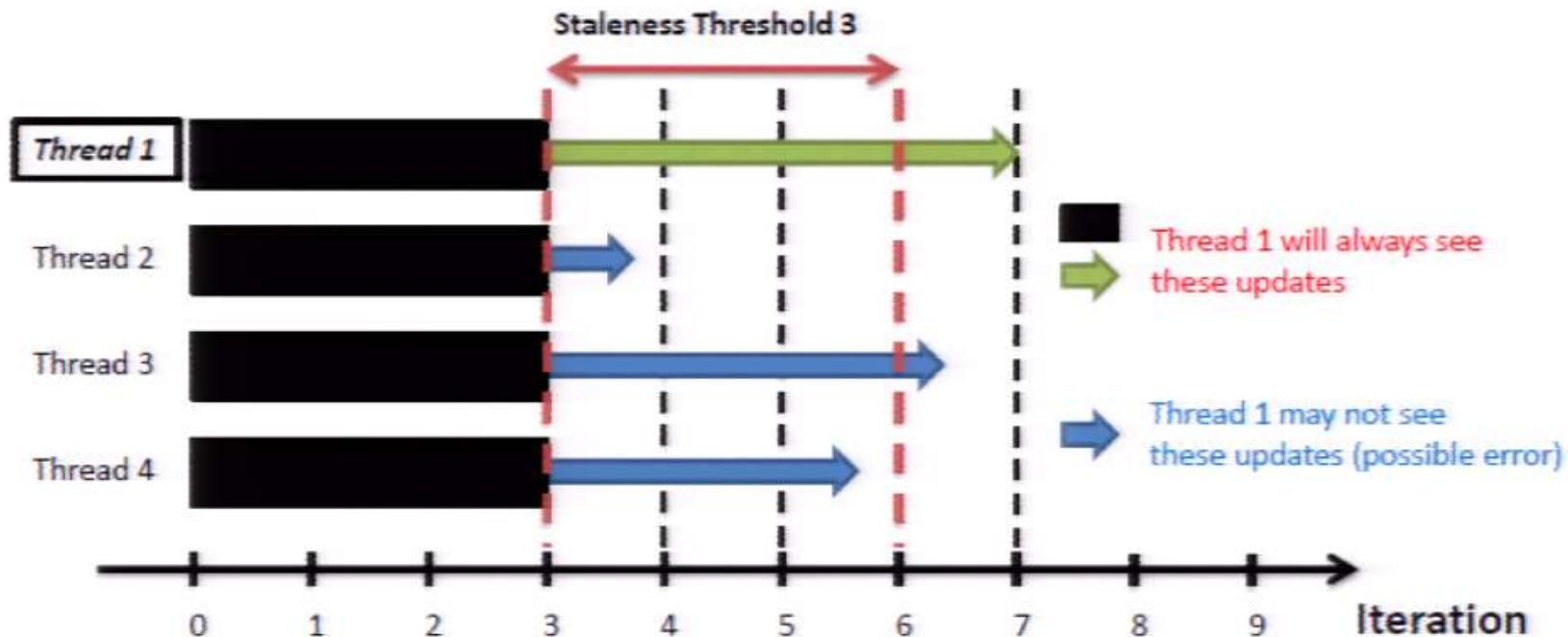
Note: x-axis is now iteration count, not time!

Allow threads to usually run at own pace

Fastest/slowest threads not allowed to drift >S iterations apart

Threads cache local (stale) versions of the parameters, to reduce network syncing

Stale Synchronous Parallel

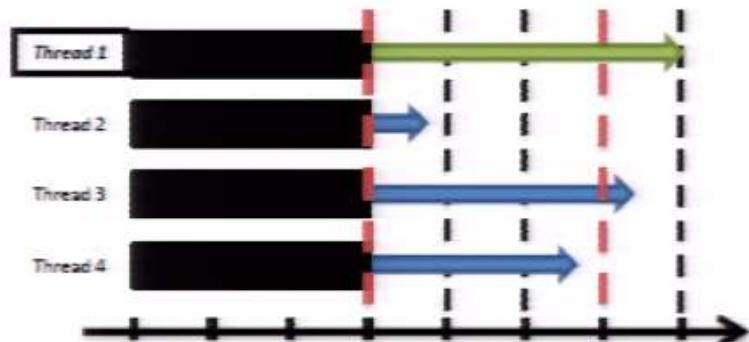


A thread at iter T sees all parameter updates before iter T-S
 Protocol: check cache first; if too old, get latest version from network

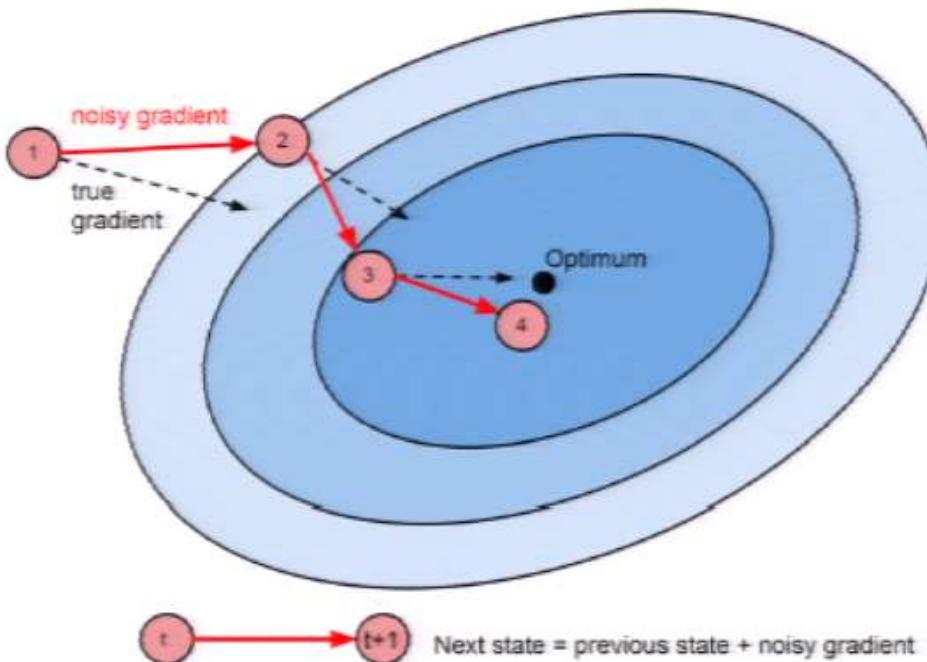
Consequence: fast threads must check network every iteration
Slow threads only check every S iterations – fewer network accesses, so catch up!

SSP provides best-of-both-worlds

- SSP combines best properties of BSP and Async
- BSP-like convergence guarantees
 - Threads cannot drift more than S iterations apart
 - Every thread sees all updates before iteration $T-S$
- Asynchronous-like speed
 - Threads usually don't wait (unless there is drift)
 - Slower threads read from network less often, thus catching up
- SSP is a spectrum of choices
 - Can be fully synchronous ($S = 0$) or very asynchronous ($S \rightarrow \infty$)
 - Or just take the middle ground, and benefit from both!



Why does SSP converge?

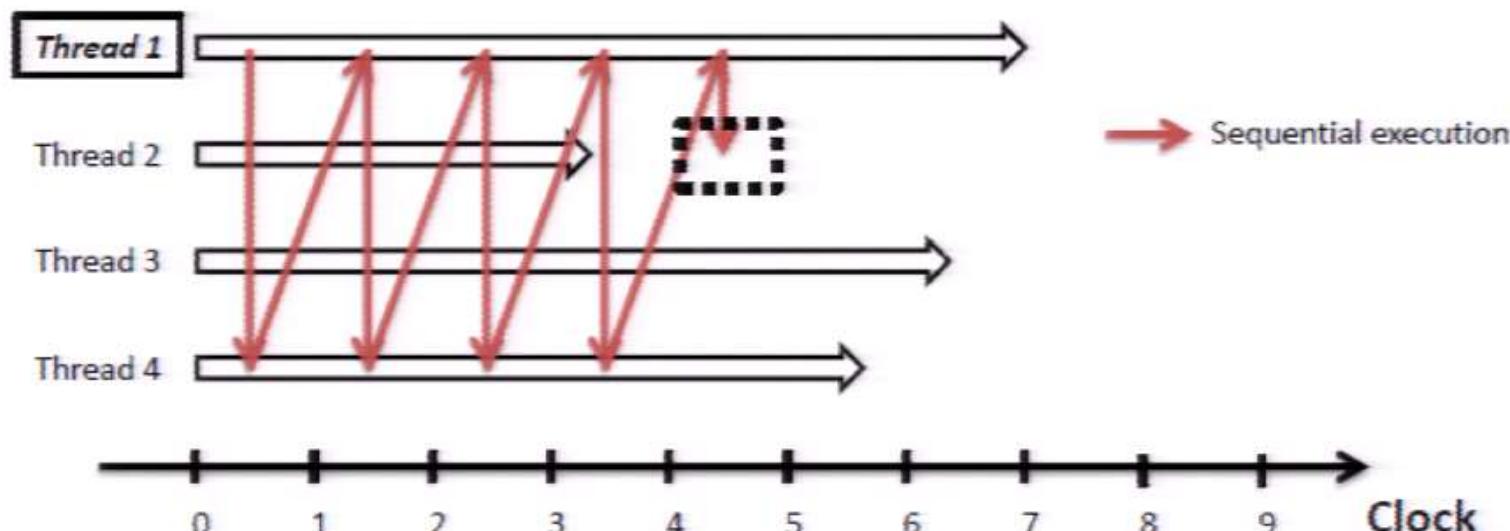


Instead of x_{true} , SSP sees $x_{stale} = x_{true} + \text{error}$

The *error caused by staleness is bounded*
Over many iterations, average error goes to zero

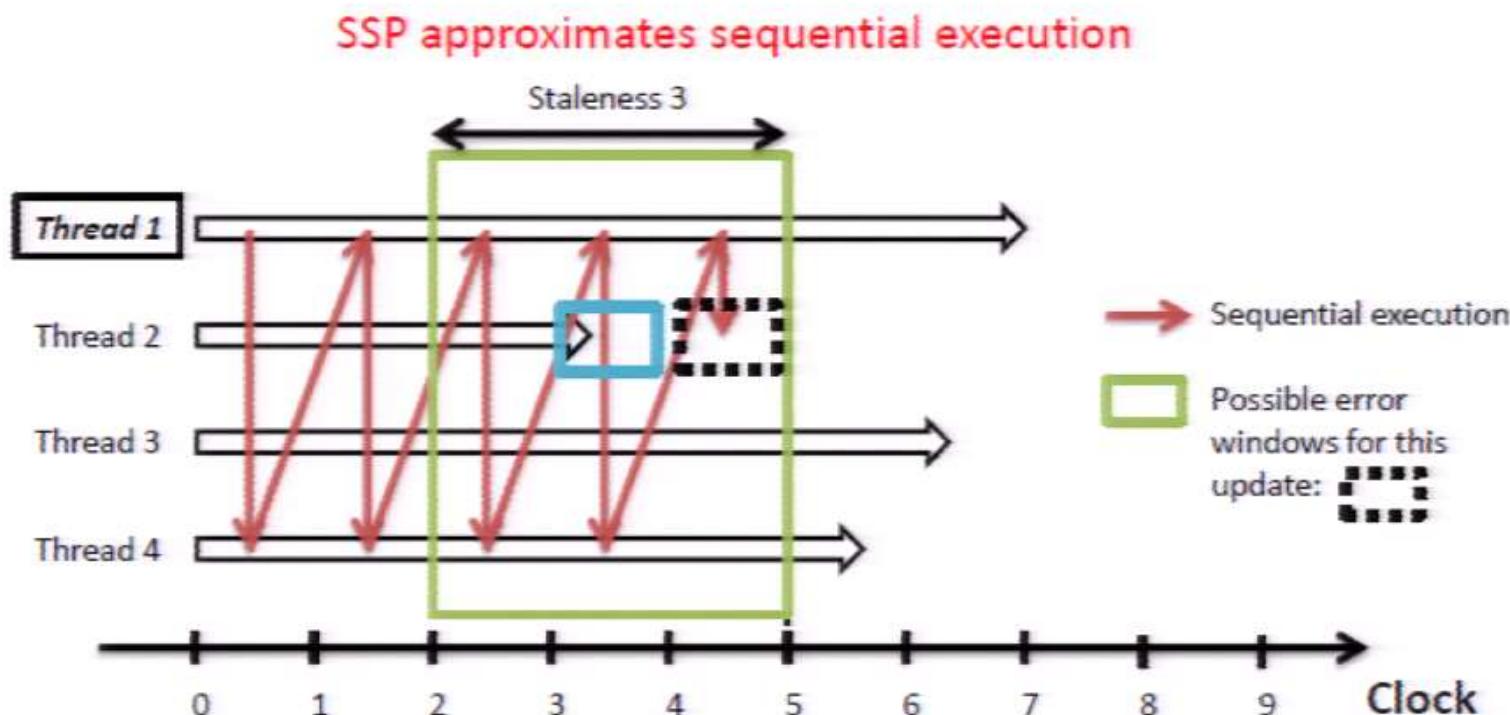
Why does SSP converge?

SSP approximates sequential execution



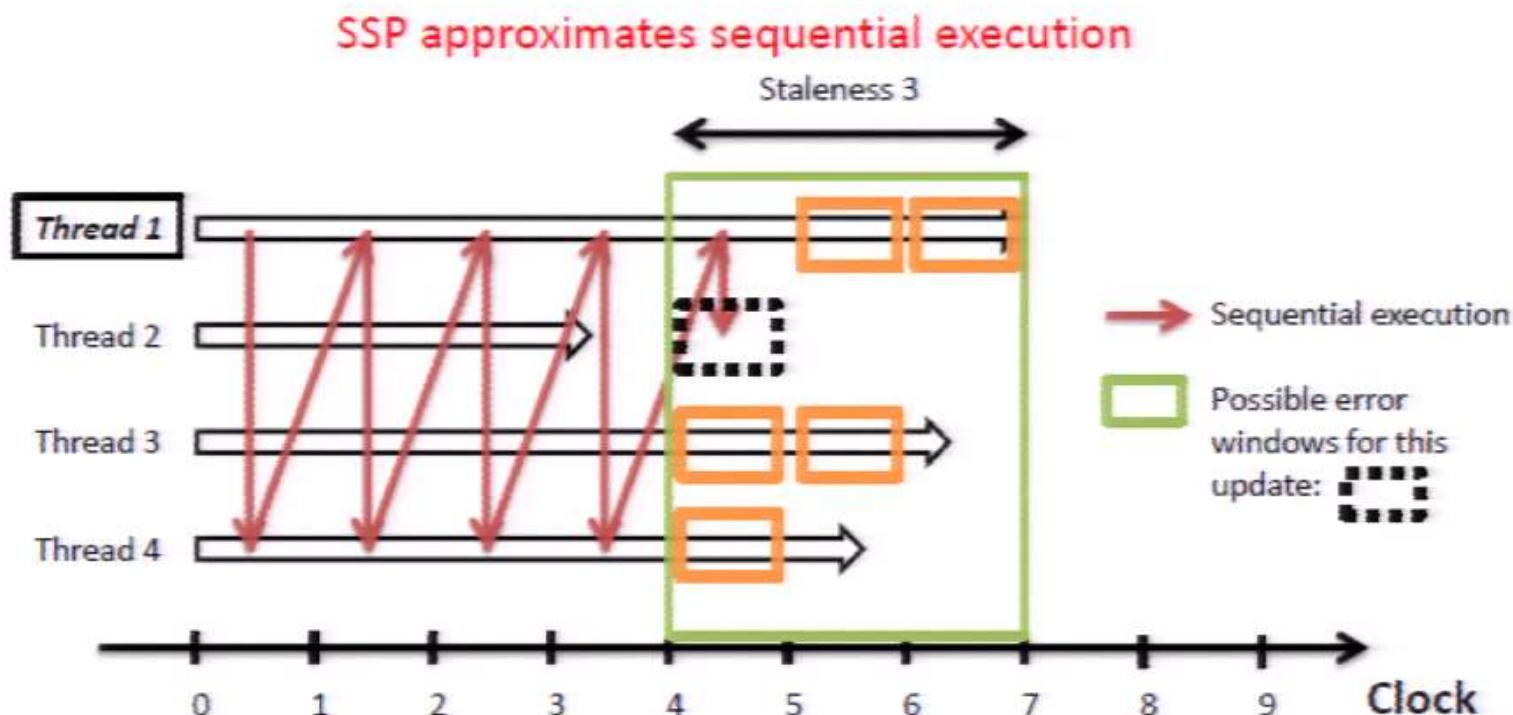
Compare actual update order to ideal sequential execution

Why does SSP converge?



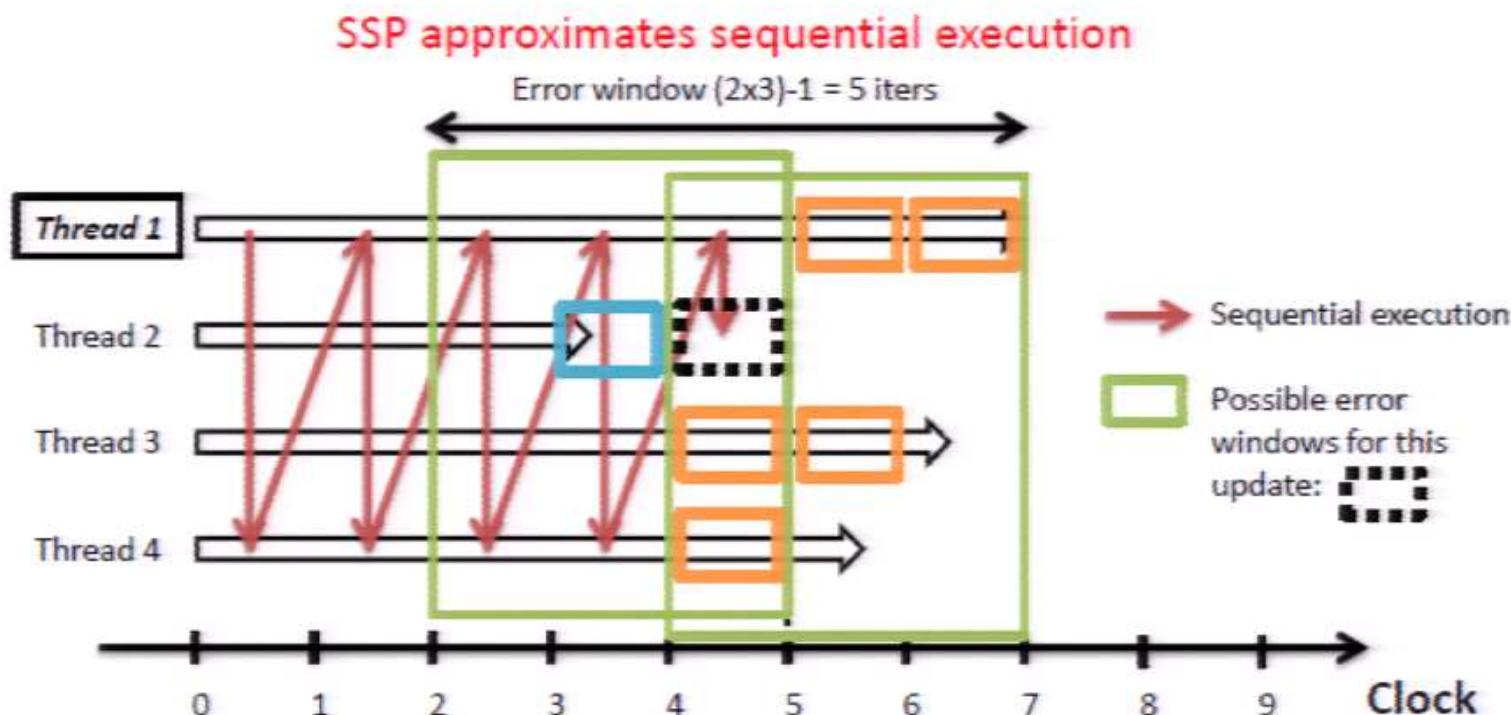
SSP may lose up to S iterations of updates to the left...

Why does SSP converge?



... as well as gain up to 5 iterations of updates to the right

Why does SSP converge?



Thus, at most $2S-1$ iterations of erroneous updates
 Hence numeric error in parameters is also bounded
 Partial, but bounded, loss of serializability

Convergence Theorem

- Want: minimize convex $f(\mathbf{x}) = \frac{1}{T} \sum_{t=1}^T f_t(\mathbf{x})$ (Example: Stochastic Gradient)
 - L -Lipschitz, problem diameter bounded by F^2
 - Staleness s , using P threads across all machines
 - Use step size $\eta_t = \frac{\sigma}{\sqrt{t}}$ with $\sigma = \frac{F}{L\sqrt{2(s+1)P}}$

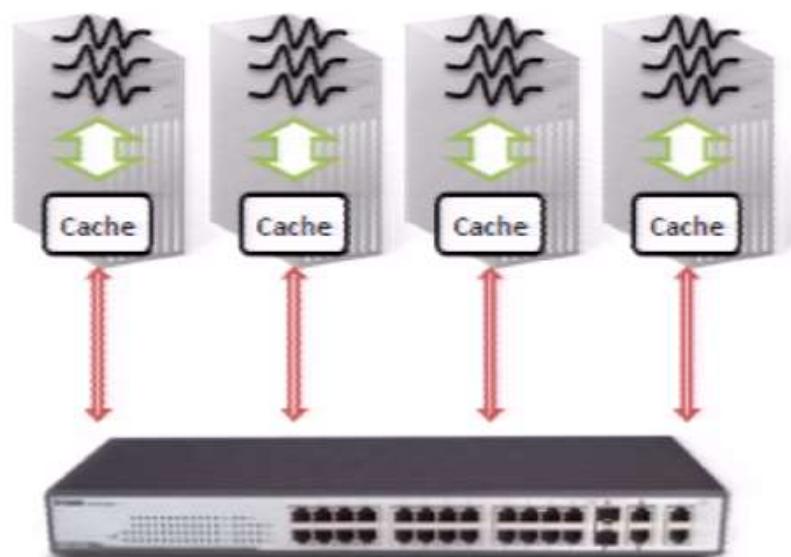
Convergence Theorem

- Want: minimize convex $f(\mathbf{x}) = \frac{1}{T} \sum_{t=1}^T f_t(\mathbf{x})$ (Example: Stochastic Gradient)
 - L -Lipschitz, problem diameter bounded by F
 - Staleness s , using P threads across all machines
 - Use step size $\eta_t = \frac{\sigma}{\sqrt{t}}$ with $\sigma = \frac{F}{L\sqrt{2(s+1)P}}$

- Difference between
SSP estimate and true optimum
- $$\overbrace{\quad\quad\quad}^{R[\mathbf{X}]} := \left[\frac{1}{T} \sum_{t=1}^T f_t(\hat{\mathbf{x}}_t) \right] - f(\mathbf{x}^*) \leq 4FL\sqrt{\frac{2(s+1)P}{T}}$$
- SSP converges according to
 - Where T is the number of iterations
 - Note: RHS bound contains (L, F) and (s, P)
 - The interaction between theory and systems parameters

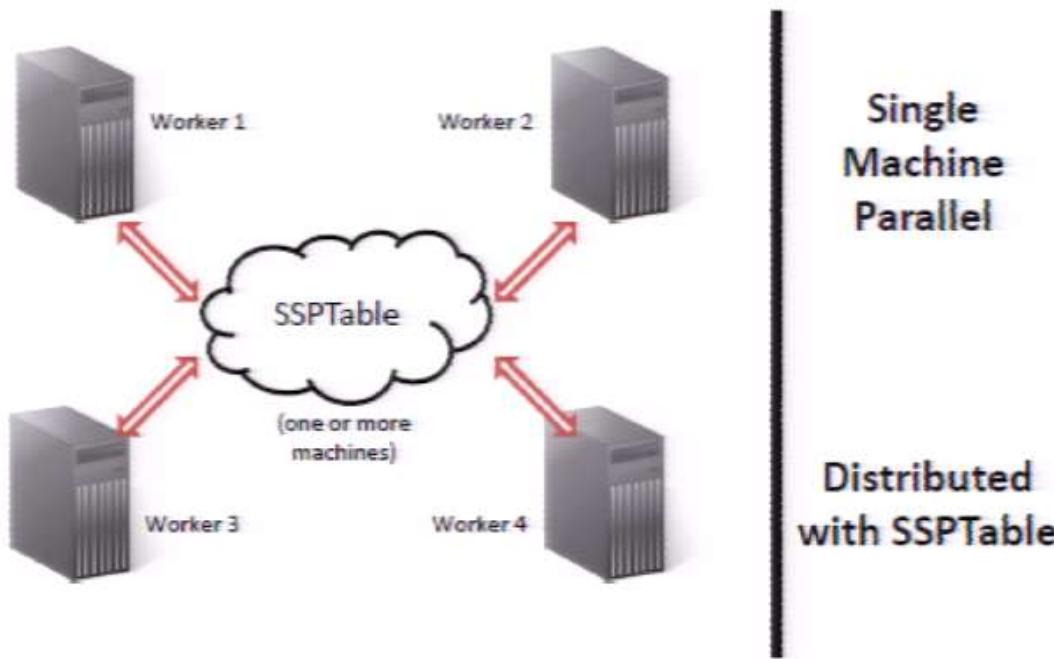
SSP solves Distributed ML challenges

- SSP is a synchronization model for fast and correct distributed ML
 - For “abelian” parameter updates of the form $\theta_{\text{new}} = \theta_{\text{old}} + \Delta$
- SSP reduces network traffic
 - Threads use stale local cache whenever possible
 - Addresses slow network and occasional machine slowdowns



SSP + Parameter Server = Easy Distributed ML

- We implement SSP as a “parameter server” (PS), called **SSPTable**
 - Provides all machines with convenient access to global model parameters
 - Can be run on multiple machines – reduces load per machine
- SSPTable allows easy conversion of single-machine parallel ML algorithms**
 - “Distributed shared memory” programming style
 - No need for complicated message passing
 - Replace local memory access with PS access



Single
Machine
Parallel

```
UpdateVar(i) {
    old = y[i]
    delta = f(old)
    y[i] += delta
}
```

Distributed
with SSPTable

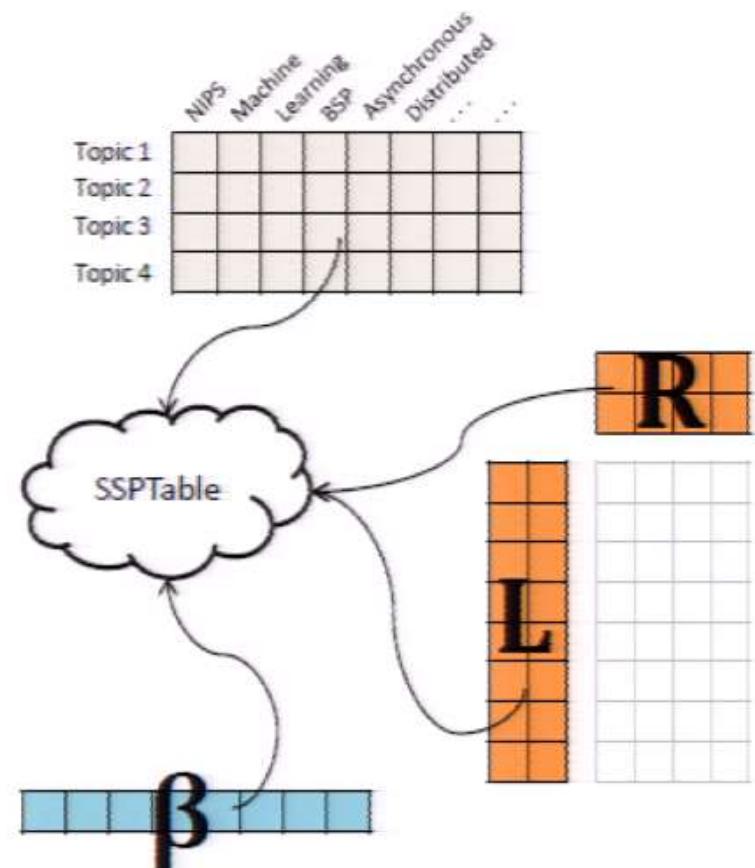
```
UpdateVar(i) {
    old = PS.read(y,i)
    delta = f(old)
    PS.inc(y,i,delta)
}
```

SSPTable Programming

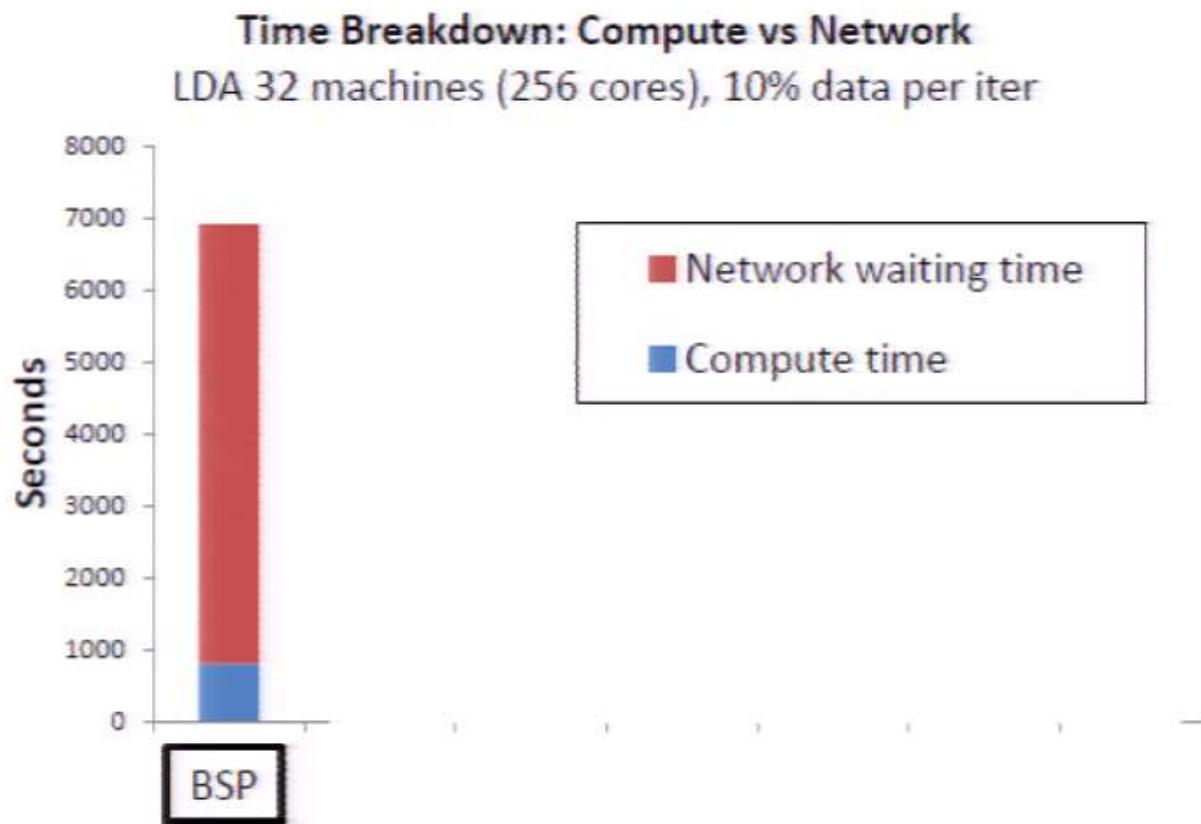
- Easy, table-based programming – just 3 commands!
 - No message passing, barriers, locks, etc.
- `read_row(table, row, s)`
 - Retrieve a table row with staleness s
- `inc(table, row, col, value)`
 - Increment table's (row,col) by value
- `clock()`
 - Inform PS that this thread is advancing to the next iteration

SSPTable Programming

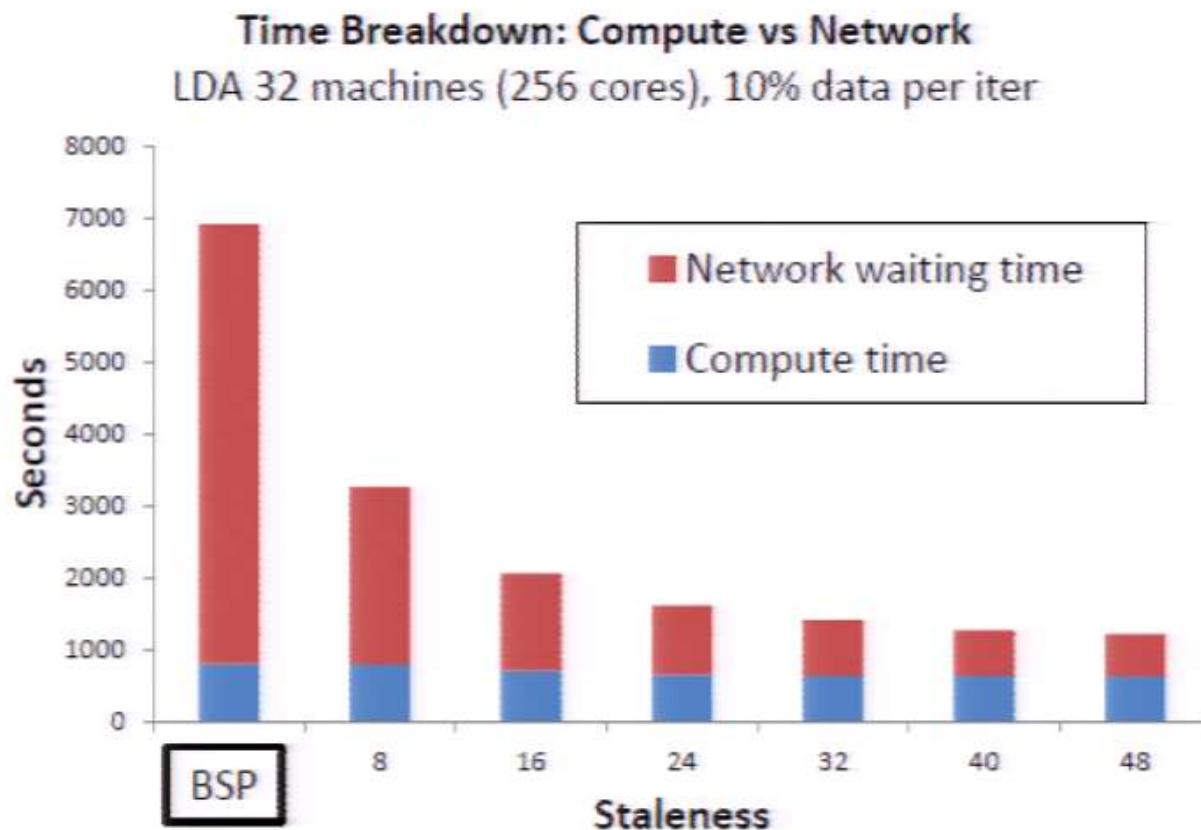
- Just put global parameters in SSPTable! Examples:
- **Topic Modeling (MCMC)**
 - Topic-word table
- **Matrix Factorization (SGD)**
 - Factor matrices L, R
- **Lasso Regression (CD)**
 - Coefficients β
- SSPTable supports **generic classes** of algorithms
 - With these models as examples



SSPTable uses networks efficiently

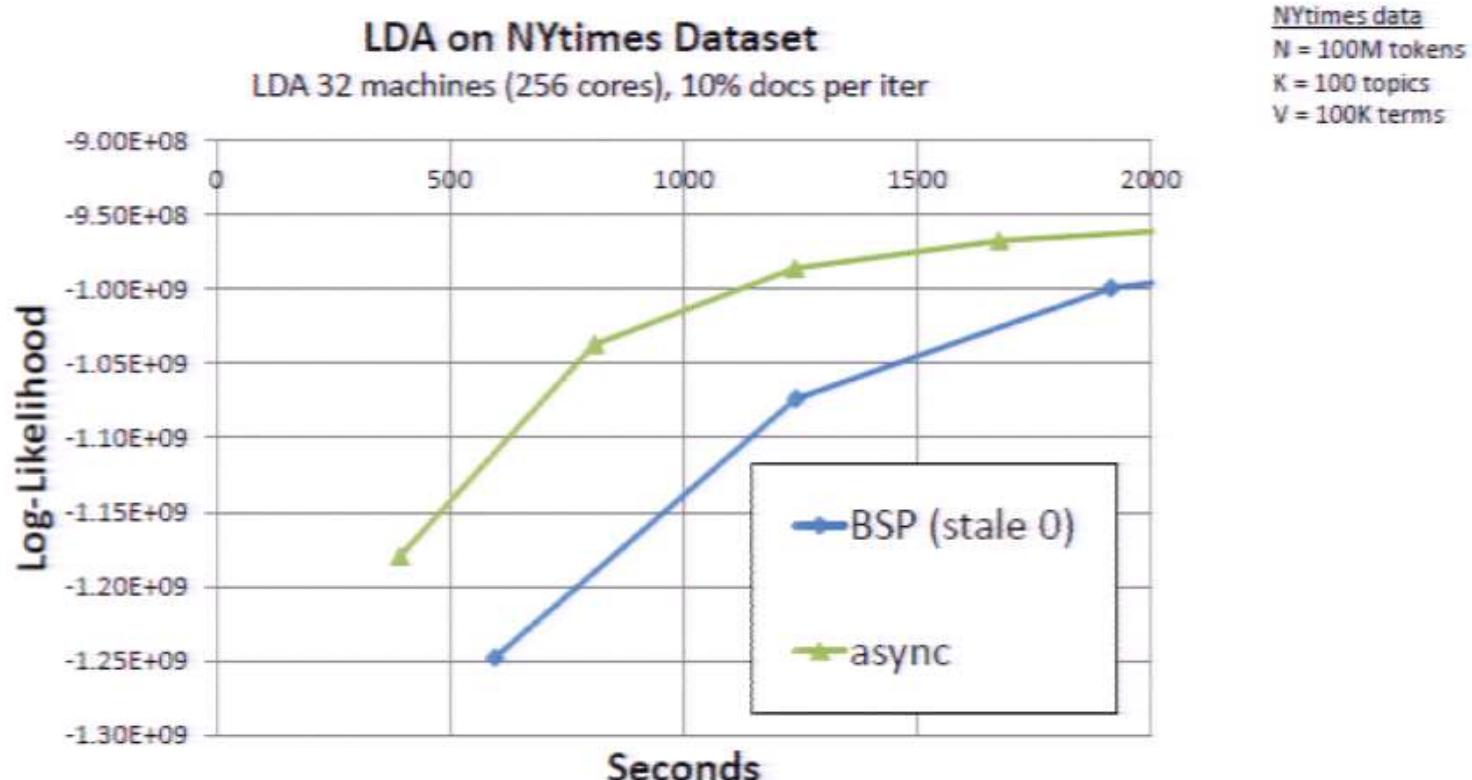


SSPTable uses networks efficiently



Network communication is a huge bottleneck with many machines
SSP balances network and compute time

SSPTable vs BSP and Async

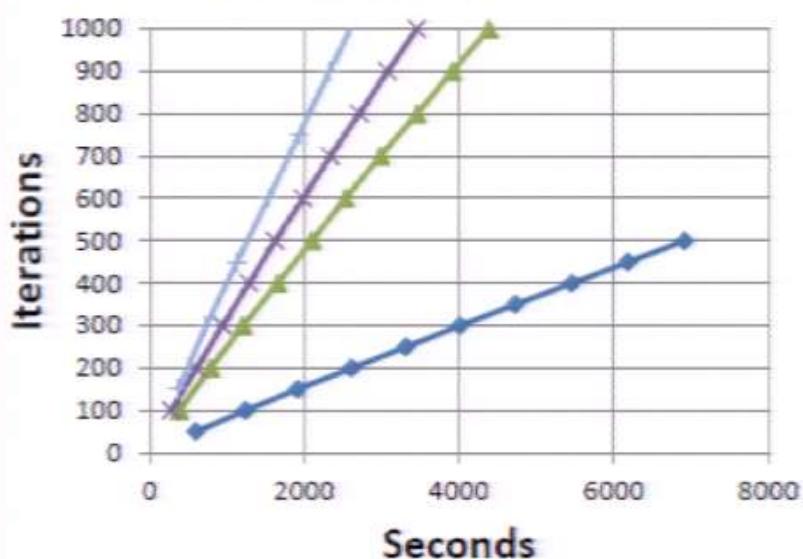


BSP has strong convergence guarantees but is slow
Asynchronous is fast but has weak convergence guarantees

The Quality vs Quantity tradeoff

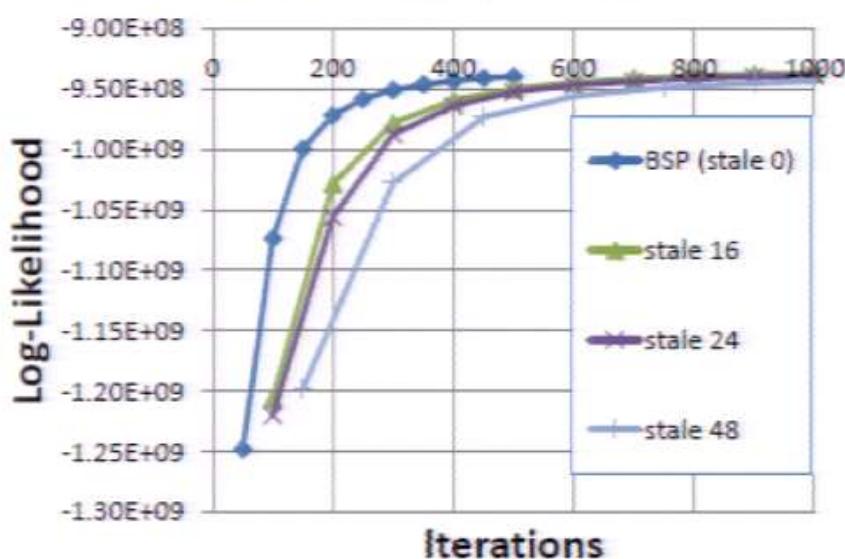
Quantity: iterations versus time

LDA 32 machines, 10% data



Quality: objective versus iterations

LDA 32 machines, 10% data

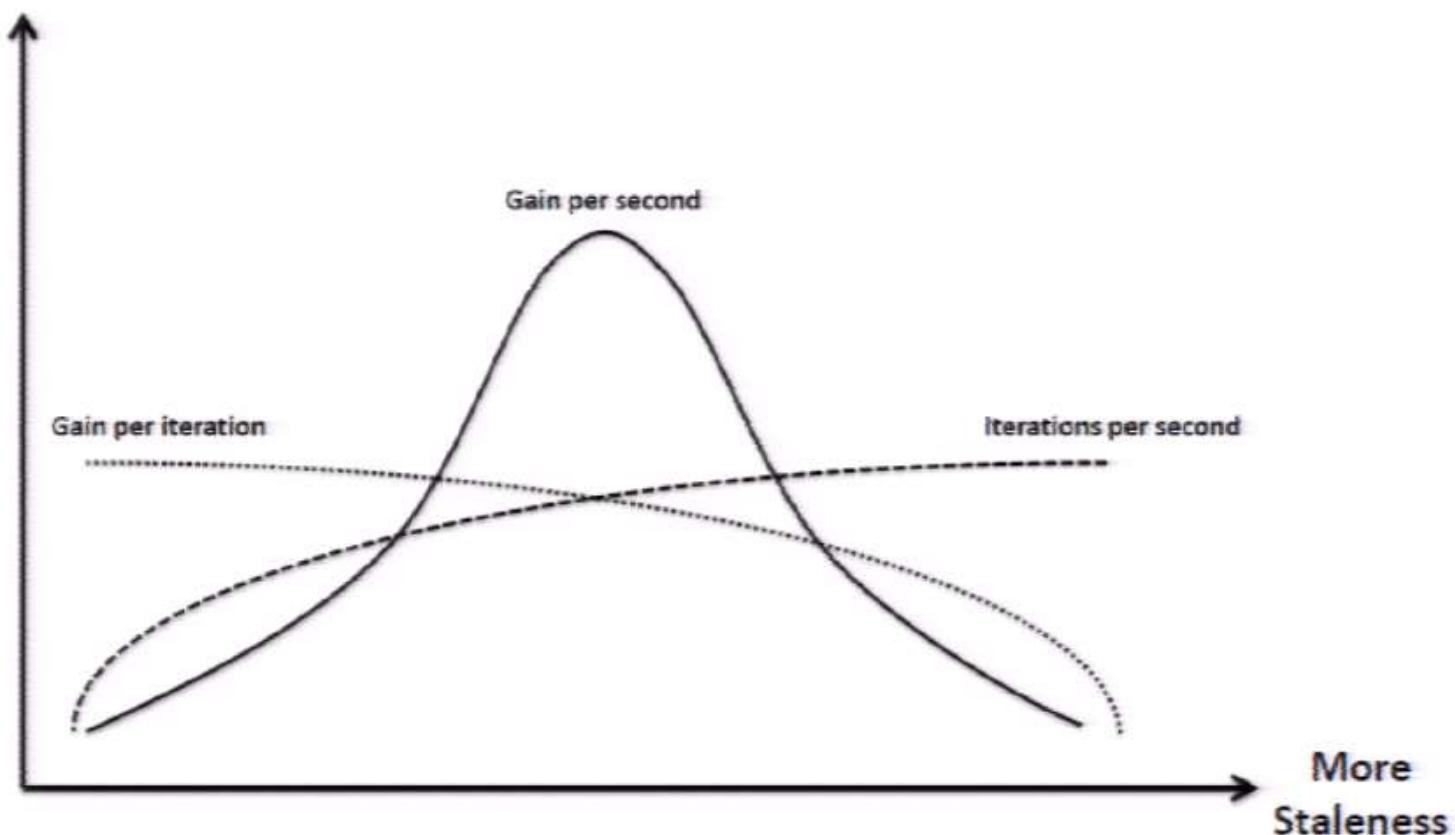


Progress per time is (iters/sec) * (progress/iter)

High staleness yields more iters/sec, but lowers progress/iter

Find the sweet spot staleness >0 for maximum progress per second

The Quality vs Quantity tradeoff



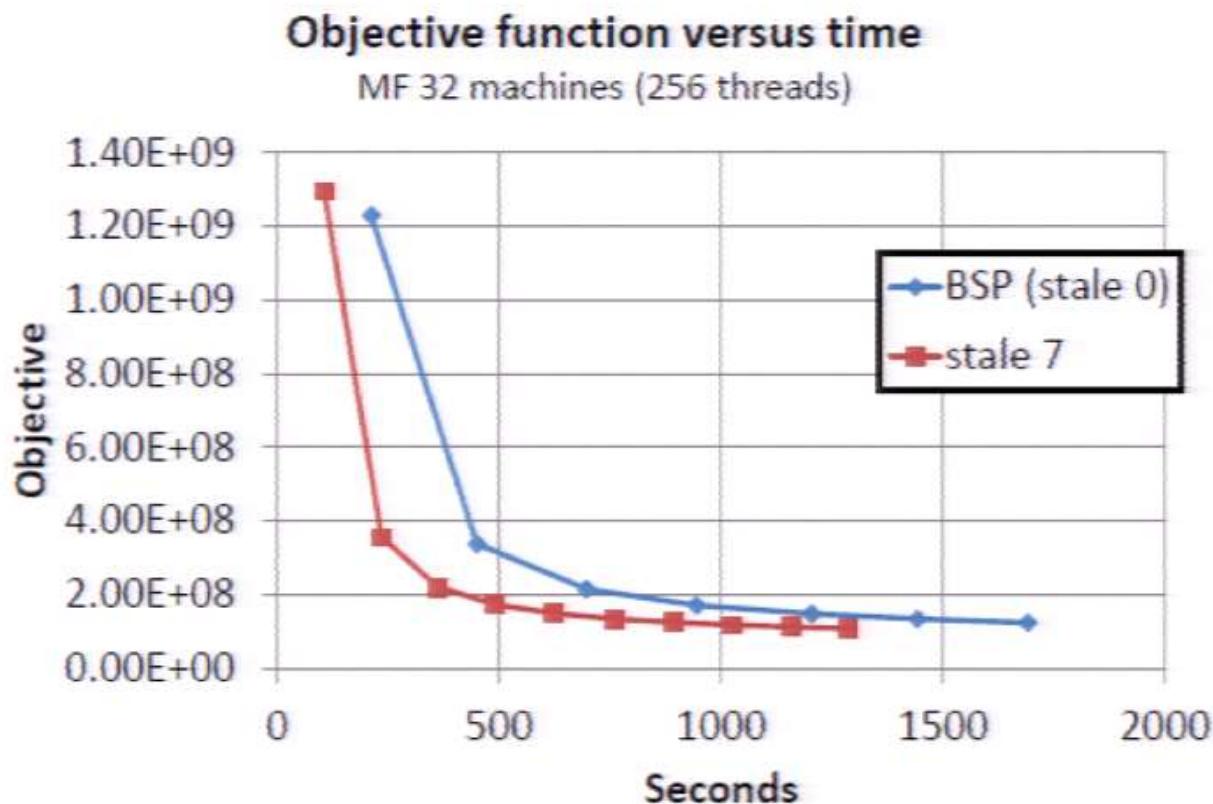
Progress per time is $(\text{iters/sec}) * (\text{progress/iter})$

High staleness yields more iters/sec, but lowers progress/iter

Find the sweet spot staleness >0 for maximum progress per second

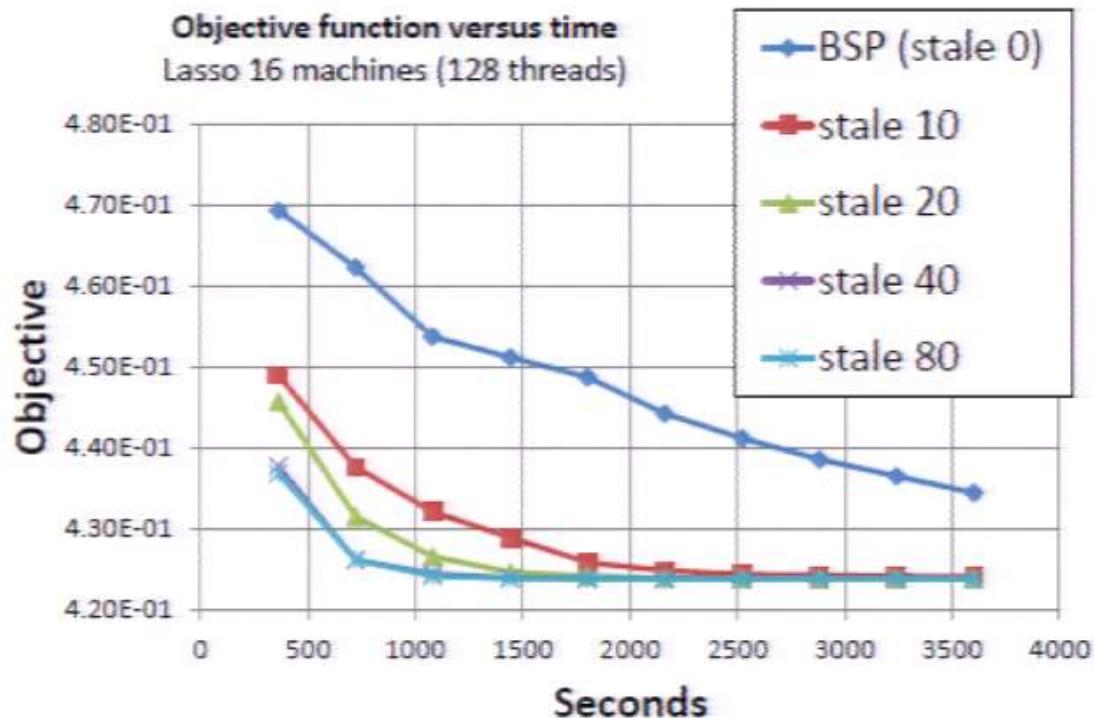
Matrix Factorization (Netflix)

Netflix data
 100M nonzeros
 480K rows
 18K columns
 rank 100



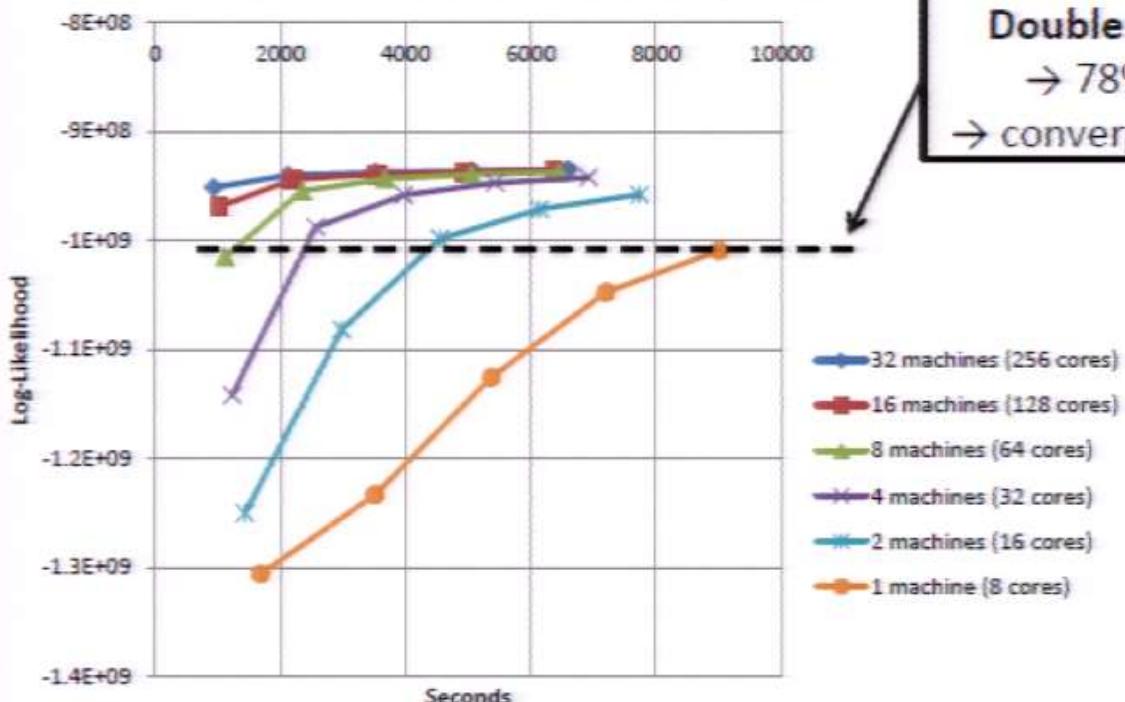
LASSO (Synthetic)

Synthetic data
 $N = 500$ samples
 $P = 400K$ features

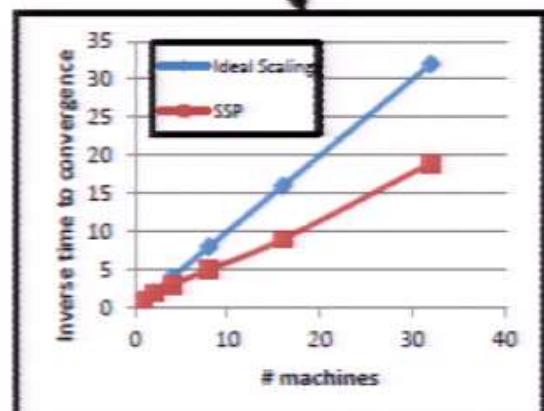


SSPTable scaling with # machines

LDA on NYtimes dataset
(staleness = 10, 1k docs per core per iteration)



Double # machines:
→ 78% speedup
→ converge in 56% time



SSP computational model scales with increasing # machines
(given a fixed dataset)

Recent Results

- Using 8 machines * 16 cores = 128 threads
 - 128GB RAM per machine
- Latent Dirichlet Allocation
 - NYTimes dataset (100M tokens, 100K words, 10K topics)
 - SSP 100K tokens/s
 - GraphLab 80K tokens/s
 - PubMed dataset (7.5B tokens, 141K words, 100 topics)
 - SSP 3.3M tokens/s
 - GraphLab 1.8M tokens/s
- Network latent space role modeling
 - Friendster network sample (39M nodes, 180M edges)
 - 50 roles: SSP takes 14h to converge (vs 5 days on one machine)

Future work

- Theory
 - SSP for MCMC
 - Automatic staleness tuning
 - Average-case analysis for better bounds
- Systems
 - Load balancing
 - Fault tolerance
 - Prefetching
 - Other consistency schemes
- Applications
 - Hard-to-parallelize ML models
 - DNNs, Regularized Bayes, Network Analysis models

Coauthors



James Cipar



Henggang Cui



Jin Kyu Kim



Seunghak Lee



Phillip B.
Gibbons



Garth A. Gibson



Gregory R.
Ganger



Eric P. Xing



Workshop Demo

- SSP is part of a bigger system: Petuum
 - SSP parameter server
 - STRADS dynamic variable scheduler
 - More features in the works
- We have a demo!
 - Topic modeling (8.2M docs, 7.5B tokens, 141K words, 10K topics)
 - Lasso regression (100K samples, 100M dimensions, 5 billion nonzeros)
 - Network latent space modeling (39M nodes, 180M edges, 50 roles)
- At BigLearning 2013 workshop (Monday)
 - <http://biglearn.org/>

Summary

- **Distributed ML is nontrivial**
 - Slow network
 - Unequal machine performance
- **SSP addresses those problems**
 - Efficiently use network resources; reduces waiting time
 - Allows slow machines to catch up
 - Fast like Async, converges like BSP
- **SSPTable parameter server provides easy table interface**
 - Quickly convert single-machine parallel ML algorithms to distributed
- Slides: www.cs.cmu.edu/~qho/ssp_nips2013.pdf

NIPS Thanks Its Sponsors



amazon.com

Microsoft
Research

Google

facebook

SKYTREE.
THE MACHINE LEARNING COMPANY

TWO²SIGMA

United Technologies
Research Center

YAHOO!
LABS

IBM
Research

xerox

DE Shaw & Co

millionshort



DRW TRADING GROUP

criteo

PDT PARTNERS

Springer
Machine Learning Journal

Mickey Mouse icon
Disney Research