

# Primitives for Dynamic Big Model Parallelism

Seunghak Lee, Jin Kyu Kim, Xun Zheng, Qirong Ho, Garth A. Gibson, Eric P. Xing\*

School of Computer Science

Carnegie Mellon University, Pittsburgh, PA, U.S.A.

\*email: [epxing@cs.cmu.edu](mailto:epxing@cs.cmu.edu)

June 19, 2014

## Abstract

When training large machine learning models with many variables or parameters, a single machine is often inadequate since the model may be too large to fit in memory, while training can take a long time even with stochastic updates. A natural recourse is to turn to distributed cluster computing, in order to harness additional memory and processors. However, naive, unstructured parallelization of ML algorithms can make inefficient use of distributed memory, while failing to obtain proportional convergence speedups — or can even result in divergence. We develop a framework of primitives for dynamic model-parallelism, STRADS, in order to explore partitioning and update scheduling of model variables in distributed ML algorithms — thus improving their memory efficiency while presenting new opportunities to speed up convergence without compromising inference correctness. We demonstrate the efficacy of model-parallel algorithms implemented in STRADS versus popular implementations for Topic Modeling, Matrix Factorization and Lasso.

## 1. INTRODUCTION

Sensory techniques and digital storage media have improved at a breakneck pace, leading to massive “Big Data” collections that have been the focus of recent efforts to achieve scalable machine learning (ML). Numerous *data-parallel* algorithmic and system solutions, both heuristic and principled, have been proposed to speed up inference on Big Data [6, 14, 16, 22]; however, large-scale ML also encompasses *Big Model* problems [7], in which models with millions if not billions of variables and/or parameters (such as in deep networks [5] or large-scale topic models [17]) must be estimated from big (or even modestly-sized) datasets. These Big Model problems seem to have received less attention in ML communities, which, in turn, has limited their application to real-world problems.

Big Model problems are challenging because a large number of model variables must be efficiently updated until model convergence. Data-parallel algorithms such as stochastic gradient descent [24] concurrently update all model variables given a subset of data samples, but this requires every worker to have full access to all global variables — which can be very large, such as the billions of variables in Deep Neural Networks [5], or this paper’s large scale topic model with 22M bigrams by 10K topics (200 *billion* variables) and matrix factorization with rank 2K on a 480K-by-10K matrix (1B variables). Furthermore, data-parallelism does not consider the possibility that some variables may be more important than others for algorithm convergence, a point that we shall demonstrate through our Lasso implementation (run on 100M coefficients). On the other hand, model-parallel algorithms such as coordinate descent [4] are well-suited to Big Model problems, because parallel workers focus on subsets of model variables. This allows the variable space to be partitioned for memory efficiency, and also allows some variables to be prioritized over others. However, model-parallel algorithms are usually developed for a specific application such as Matrix Factorization [9] or Lasso [4] — thus, there is utility in developing programming primitives that can tackle the common challenges of Big Model problems, while also exposing new opportunities such as variable prioritization.

Existing distributed frameworks such as MapReduce [6] and GraphLab [14] have shown that

Table 1: Summary of LDA, MF, and Lasso on STRADS (detailed pseudocode is in the relevant sections).

	<b>Schedule</b>	<b>Push and Pull</b>	<b>Largest STRADS experiment</b>
Topic Modeling (LDA)	Word rotation scheduling	Collapsed Gibbs sampling	10K topics, 3.9M docs with 21.8M vocab
MF	Round-robin scheduling	Coordinate descent	rank-2K, 480K-by-10K matrix
Lasso	Dynamic priority scheduling	Coordinate descent	100M features, 50K samples

common primitives such as Map/Reduce or Gather/Apply/Scatter can be applied to a variety of ML applications. Crucially, these frameworks automatically decide which variable to update next — MapReduce executes all Mappers at the same time, followed by all Reducers, while GraphLab chooses the next node based on its “chromatic engine” and the user’s choice of graph consistency model. While such *automatic scheduling* is convenient, it does not offer the fine-grained control needed to avoid parallelization of variables with subtle interdependencies not seen in the superficial problem or graph structure (which can then lead to algorithm divergence, as in Lasso [4]). Moreover, it does not allow users to explicitly prioritize variables based on new criteria.

To improve upon these frameworks, we develop new primitives for dynamic Big Model parallelism: **schedule**, **push** and **pull**, which are executed by our STRADS system (STRucture-Aware Dynamic Scheduler). These primitives are inspired by the simplicity and wide applicability of MapReduce, but also provide the fine control needed to explore novel ways of performing dynamic model-parallelism. **Schedule** specifies the next subset of model variables to be updated in parallel, **push** specifies how individual workers compute partial results on those variables, and **pull** specifies how those partial results are aggregated to perform the full variable update. A final “automatic primitive”, **sync**, ensures that distributed workers have up-to-date values of the model variables, and is automatically executed at the end of **pull**; the user does not need to implement **sync**. To explore the utility of STRADS, we implement **schedule**, **push** and **pull** for three popular ML applications (Table 1): Topic Modeling (LDA), Lasso, and Matrix Factorization (MF). Our goal is not to best specialized implementations in performance, but to demonstrate that STRADS primitives enable Big Model problems to be solved with modest programming effort. In particular, we tackle topic modeling with 3.9M docs, 10K topics and 21.8M vocabulary (200B variables),

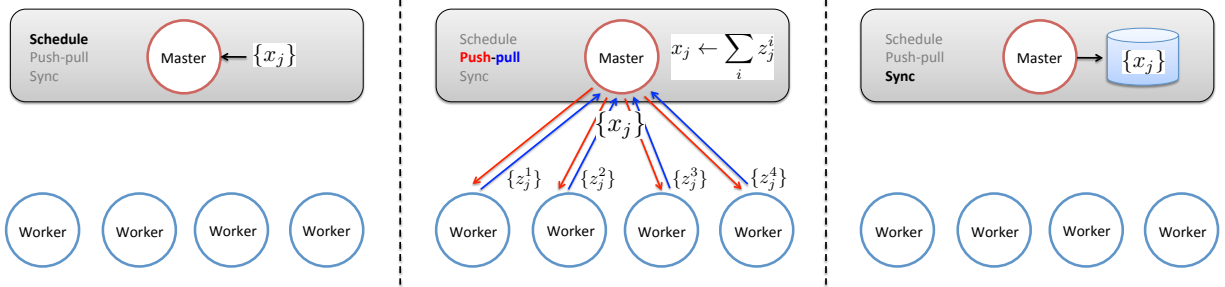


Figure 1: High-level view of our STRADS primitives for dynamic model parallelism.

MF with rank-2K on a 480K-by-10K matrix (1B variables), and Lasso with 100M features (100M variables).

## 2. PRIMITIVES FOR DYNAMIC MODEL PARALLELISM

“Model parallelism” refers to parallelization of an ML algorithm over the space of shared model variables, rather than the space of (usually i.i.d.) data samples. At a high level, model variables are the changing intermediate quantities that an ML algorithm iteratively updates, until convergence is reached. For example, the coefficients in regression are model variables, which are iteratively updated using algorithmic strategies like coordinate descent.

Model parallelism can be contrasted with data parallelism, in which the ML algorithm is parallelized over individual data samples, such as in stochastic optimization algorithms [25]. A key advantage of the model-parallel approach is that it explicitly partitions the model variables into subsets, allowing ML problems with massive model spaces to be tackled on machines with limited memory. Figure 3 shows this advantage: for topic modeling, STRADS uses less memory per machine as the number of machines increases, unlike the data-parallel YahooLDA algorithm. As our experiments will confirm, this means that STRADS can handle larger ML models (given sufficient machines), whereas YahooLDA is strictly constrained by the memory of the smallest machine. This has practical consequences — STRADS LDA can handle bigram vocabularies with over 20 million term-pairs on modest hardware (enabling large-scale topic modeling applications), while YahooLDA cannot.

```
// Generic STRADS application
```

```
schedule() {  
  // Select U vars x[j] to be sent  
  // to the workers for updating  
  ...  
  return (x[j_1], ..., x[j_U])  
}
```

```
push(worker = p, vars = (x[j_1],...,x[j_U])) {  
  // Compute partial update z for U vars x[j]  
  // at worker p  
  ...  
  return z  
}
```

```
pull(workers = [p], vars = (x[j_1],...,x[j_U]),  
      updates = [z]) {  
  // Use partial updates z from workers p to  
  // update U vars x[j]. sync() is automatic.  
  ...  
}
```

Figure 2: **STRADS user-defined primitives: schedule, push, pull.** We show the basic functional signature of each primitive, using pseudocode.

To enable users to systematically and programmatically exploit model parallelism, our proposed STRADS framework defines a set of primitives. Similar to the map-reduce paradigm, these primitives are functions that a user writes for his/her ML problem, and STRADS repeatedly executes these functions to create an iterative model-parallel algorithm (Figures 1, 2). Our primitives are **schedule**, **push** and **pull**, and a single “round” or iteration of STRADS executes them in that order. In addition, there is an automatic primitive, **sync**, which the user does not have to write.

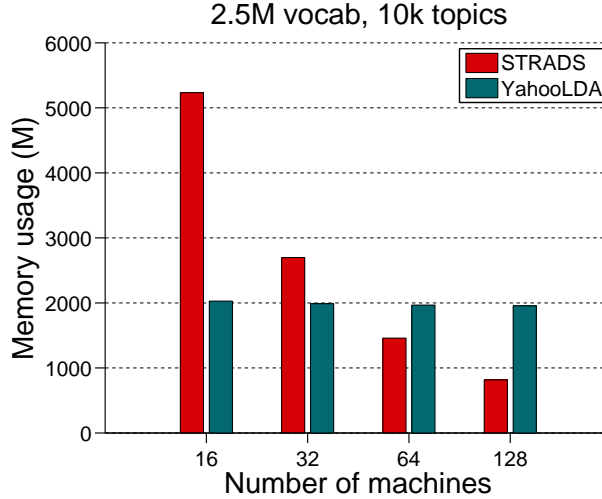


Figure 3: **Topic Modeling: Memory usage per machine**, for model-parallelism (STRADS) vs data-parallelism (YahooLDA). With more machines, STRADS LDA uses *less memory per machine*, because it explicitly partitions the model space.

**Schedule:** This primitive determines the parallel order for updating model variables; as shown in Figure 2, **schedule** selects  $U$  model variables to be dispatched for updates (Figure 1). Within the **schedule** function, the programmer may access all data  $D$  and all model variables  $x$ , in order to decide which  $U$  variables to dispatch. The simplest possible **schedule** is to select model variables according to a fixed sequence, or drawn uniformly at random. As we shall later see, **schedule** also allows model variables to be selected in a way that: (1) dynamically focuses workers on the fastest-converging variables, while avoiding already-converged variables; (2) avoids parallel dispatch of variables with inter-dependencies, which can lead to divergence and incorrect execution.

**Push and Pull:** These primitives control the flow of model variables  $x$  and data  $D$  from the master scheduler machines(s) to and from the workers (Figure 1). The **push** primitive dispatches a set of variables  $\{x_{j_1}, \dots, x_{j_U}\}$  to each worker  $p$ , which then computes a partial update  $z$  for  $\{x_{j_1}, \dots, x_{j_U}\}$  (or a subset of it). When writing **push**, the user can take advantage of data partitioning: e.g., when only a fraction  $\frac{1}{P}$  of the data samples are stored at each worker, the  $p$ -th worker should compute partial results  $z_j^p = \sum_{D_i} f_{x_j}(D_i)$  by iterating over its  $\frac{1}{L}$  data points  $D_i$ .

The **pull** primitive is used to aggregate the partial results  $\{z_j^p\}$  from all workers, and commit them to the variables  $\{x_{j_1}, \dots, x_{j_U}\}$ . Our STRADS LDA, Lasso and MF applications partition the data samples uniformly over machines.

**Synchronization:** The model variables  $x$  are globally accessible through a distributed, partitioned key-value store (represented by standard arrays in our pseudocode). **Sync** is a built-in primitive that ensures all **push** workers can access up-to-date model variables, and is automatically executed whenever **pull** writes to any variable  $x[j]$ . The user does not need to implement **sync**. A variety of key-value store synchronization schemes exist, such as Bulk Synchronous Parallel (BSP), Stale Synchronous Parallel (SSP) [13], and Asynchronous Parallel (AP). Each presents a different trade-off: BSP is simple and correct but easily bottlenecked by slow workers, AP is usually effective but risks algorithmic errors and divergence because it has no error guarantees, and SSP is fast and guaranteed to converge but requires more engineering work and parameter tuning. In this paper, we use BSP for **sync** throughout; we leave the use of alternative schemes like SSP or AP as future work.

### 3. HARNESSING MODEL-PARALLELISM IN ML APPLICATIONS THROUGH STRADS

In this section, we shall explore how users can apply model-parallelism to their own ML applications, using the STRADS primitives. We shall cover 3 ML application case studies, with the intent of showing that model-parallelism in STRADS can be simple and effective, yet also powerful enough to expose new and interesting opportunities for speeding up distributed ML.

#### 3.1 Latent Dirichlet Allocation (LDA)

We introduce STRADS programming through topic modeling via LDA [3]. Big LDA models provide a strong use case for model-parallelism: when thousands of topics and millions of words are used, the LDA model contains billions of global variables, and data-parallel implementations face the difficult challenge of providing access to all these variables; in contrast, model-parallelism explicitly divides up the variables, so that workers only need to access a fraction at a given time.

Formally, LDA takes a corpus of  $N$  documents as input, and outputs  $K$  topics (each topic is just a categorical distribution over all  $V$  unique words in the corpus) as well as  $N$   $K$ -dimensional topic vectors (soft assignments of topics to documents). The LDA model is

$$P(\mathbf{W} \mid \mathbf{Z}, \boldsymbol{\theta}, \boldsymbol{\beta}) = \prod_{i=1}^N \prod_{j=1}^{M_i} P(w_{ij} \mid z_{ij}, \boldsymbol{\beta}) P(z_{ij} \mid \boldsymbol{\theta}_i),$$

where (1)  $w_{ij}$  is the  $j$ -th token (word position) in the  $i$ -th document, (2)  $M_i$  is the number of tokens in document  $i$ , (3)  $z_{ij}$  is the topic assignment for  $w_{ij}$ , (4)  $\boldsymbol{\theta}_i$  is the topic vector for document  $i$ , and (5)  $\boldsymbol{\beta}$  is a matrix representing the  $K$   $V$ -dimensional topics. LDA is commonly reformulated as a “collapsed” model [12] in which  $\boldsymbol{\theta}, \boldsymbol{\beta}$  are integrated out for faster inference. Inference is performed using Gibbs sampling, where each  $z_{ij}$  is sampled in turn according to its distribution conditioned on all other variables,  $P(z_{ij} \mid \mathbf{W}, \mathbf{Z}_{-ij})$ . To perform this computation without having to iterate over all  $\mathbf{W}, \mathbf{Z}$ , sufficient statistics are kept in the form of a “doc-topic” table  $\mathbf{D}$  (analogous to  $\boldsymbol{\theta}$ ), and a “word-topic” table  $\mathbf{B}$  (analogous to  $\boldsymbol{\beta}$ ). More precisely,  $D_{ik}$  counts the number of assignments  $z_{ij} = k$  in doc  $i$ , while  $B_{vk}$  counts the number of tokens  $w_{ij} = v$  such that  $z_{ij} = k$ .

**STRADS implementation:** In order to perform model-parallelism, we first identify the model variables, and create a **schedule** strategy over them. In LDA, the assignments  $z_{ij}$  are the model variables, while  $\mathbf{D}, \mathbf{B}$  are summary statistics over the  $z_{ij}$  that are used to speed up the sampler. Our **schedule** strategy equally divides the  $V$  words into  $U$  subsets  $V_1, \dots, V_U$  (where  $U$  is the number of workers). Each worker will only process words from one subset  $V_a$  at a time. Subsequent invocations of **schedule** will “rotate” subsets amongst workers, so that every worker touches all  $U$  subsets every  $U$  invocations. For data partitioning, we divide the document tokens  $\mathbf{W}$  evenly across workers, and denote worker  $p$ ’s set of tokens by  $\mathbf{W}_{q_p}$ .

During **push**, suppose that worker  $p$  is assigned to subset  $V_a$  by **schedule**. This worker will only Gibbs sample the topic assignments  $z_{ij}$  such that (1)  $(i, j) \in \mathbf{W}_{q_p}$  and (2)  $w_{ij} \in V_a$ . In other words,  $w_{ij}$  must be assigned to worker  $p$ , and must also be a word in  $V_a$ . The latter condition is the source of model-parallelism: observe how the assignments  $z_{ij}$  are chosen for sampling based on word



```
// STRADS LDA
```

```
schedule() {  
  dispatch = [] // Empty list  
  for a=1..U // Rotation scheduling  
    idx = ((a+C-1) mod U) + 1  
    dispatch.append( V[q_idx] )  
  return dispatch  
}
```

```
push(worker = p, vars = [V_a, ..., V_U]) {  
  t = [] // Empty list  
  for (i,j) in W[q_p] // Fast Gibbs sampling  
    if w[i,j] in V_p  
      t.append( (i, j, f_1(i, j, D, B)) )  
  return t  
}
```

```
pull(workers = [p], vars = [V_a, ..., V_U],  
      updates = [t]) {  
  for all (i,j) // Update sufficient stats  
    (D,B) = f_2([t])  
}
```

Figure 4: **STRADS LDA pseudocode.** Definitions for  $f_1, f_2, q_p$  are in the text.  $C$  is a global model variable.

divisions  $V_a$ . Note that all  $z_{ij}$  will be sampled exactly once after  $U$  invocations of **schedule**. We use the fast Gibbs sampler from [20] to **push** update  $z_{ij} \leftarrow f_1(i, j, \mathbf{D}, \mathbf{B})$ , where  $f_1(\cdot)$  represents the fast Gibbs sampler equation. The **pull** step simply updates the sufficient statistics  $\mathbf{D}, \mathbf{B}$  using the new  $z_{ij}$ , and we represent this procedure as a function  $(\mathbf{D}, \mathbf{B}) \leftarrow f_2([z_{ij}])$ . Figure 4 provides pseudocode for STRADS LDA.

**Model parallelism results in low error:** Parallel Gibbs sampling is not generally guaranteed to converge [11], unless the variables being parallel-sampled are conditionally independent of each

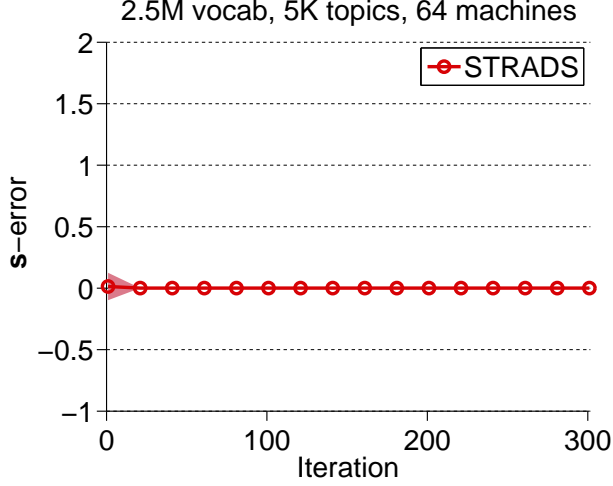


Figure 5: **STRADS LDA**:  $s$ -error  $\Delta_{r,t}$  at each iteration, on the Wikipedia unigram dataset with  $K = 5000$  and 64 machines.

other. Because STRADS LDA assigns workers to disjoint words  $V$  and documents  $w_{ij}$ , each worker’s variables  $z_{ij}$  are (almost) conditionally independent of other workers, except for a single shared dependency: the column sums of  $\mathbf{B}$  (denoted by  $\mathbf{s}$ , and stored as an extra row appended to  $\mathbf{B}$ ), which are required for correct normalization of the Gibbs sampler conditional distributions in  $f_1()$ . The column sums  $\mathbf{s}$  are **synced** at the end of every **pull**, but will go out-of-sync during worker **pushes**. To understand how error in  $\mathbf{s}$  affects sampler convergence, consider the Gibbs sampling conditional distribution for a topic indicator  $z_{ij}$ :

$$\begin{aligned} P(z_{ij} \mid \mathbf{W}, \mathbf{Z}_{-ij}) &\propto P(w_{ij} \mid z_{ij}, \mathbf{W}_{-ij}, \mathbf{Z}_{-ij})P(z_{ij} \mid \mathbf{Z}_{-ij}) \\ &= \frac{\gamma + B_{w_{ij}, z_{ij}}}{V\gamma + \sum_{v=1}^V B_{v, z_{ij}}} \times \frac{\alpha + D_{i, z_{ij}}}{K\alpha + \sum_{k=1}^K D_{i, k}}. \end{aligned}$$

In the first term, the denominator quantity  $\sum_{v=1}^V B_{v, z_{ij}}$  is exactly the sum over the  $z_{ij}$ -th column of  $\mathbf{B}$ , i.e.  $\mathbf{s}_{z_{ij}}$ . Thus, errors in  $\mathbf{s}$  induce errors in the probability distribution  $U_{w_{ij}} \sim P(w_{ij} \mid z_{ij}, \mathbf{W}_{-ij}, \mathbf{Z}_{-ij})$ , which is just the discrete probability that topic  $z_{ij}$  will generate word  $w_{ij}$ . As a proxy for the error in  $U$ , we can measure the difference between the true  $\mathbf{s}$  and its local copy  $\tilde{\mathbf{s}}^p$  on worker  $p$ . If  $\mathbf{s} = \tilde{\mathbf{s}}^p$ , then  $U$  has zero error.

We can show that the error in  $\mathbf{s}$  is empirically negligible (and hence the error in  $U$  is also small).

Consider a single STRADS LDA iteration  $t$ , and define its  $s$ -error to be

$$\Delta_t = \frac{1}{PM} \sum_{p=1}^P \|\tilde{\mathbf{s}}^p - \mathbf{s}\|_1, \quad (1)$$

where  $M$  is the total number of tokens  $w_{ij}$ . The  $s$ -error  $\Delta_{r,t}$  must lie in  $[0, 2]$ , where 0 means no error. Figure 5 plots the  $s$ -error for the ‘‘Wikipedia unigram’’ dataset (refer to our experiments section for details), for  $K = 5000$  topics and 64 machines (128 processor cores total). The  $s$ -error is  $\leq 0.002$  throughout, confirming that STRADS LDA exhibits very small parallelization error.

### 3.2 Matrix Factorization (MF)

STRADS’s model-parallelism benefits other models as well: we now consider Matrix Factorization (collaborative filtering), which can be used to predict users’ unknown preferences, given their known preferences and the preferences of others. While most MF implementations tend to focus on small decompositions with rank  $K \approx 100$  [23, 9, 21], we are interested in enabling larger decompositions with rank  $> 1000$ , where the much larger factors (billions of variables) pose a challenge for purely data-parallel algorithms (such as naive SGD) that need to share all variables across all workers; again, STRADS addresses this by explicitly dividing variables across workers.

Formally, MF takes an incomplete matrix  $\mathbf{A} \in \mathbb{R}^{N \times M}$  as input, where  $N$  is the number of users, and  $M$  is the number of items/preferences. The idea is to discover rank- $K$  matrices  $\mathbf{W} \in \mathbb{R}^{N \times K}$  and  $\mathbf{H} \in \mathbb{R}^{K \times M}$  such that  $\mathbf{WH} \approx \mathbf{A}$ . Thus, the product  $\mathbf{WH}$  can be used to predict the missing entries (user preferences). Formally, let  $\Omega$  be the set of indices of observed entries in  $\mathbf{A}$ , let  $\Omega^i$  be the set of observed column indices in the  $i$ -th row of  $\mathbf{A}$ , and let  $\Omega_j$  be the set of observed row indices in the  $j$ -th column of  $\mathbf{A}$ . Then, the MF task is defined as an optimization problem:

$$\min_{\mathbf{W}, \mathbf{H}} \sum_{(i,j) \in \Omega} (a_j^i - \mathbf{w}^i \mathbf{h}_j)^2 + \lambda (\|\mathbf{W}\|_F^2 + \|\mathbf{H}\|_F^2). \quad (2)$$

This can be solved using parallel CD [21], with the following update rule for  $\mathbf{H}$ :

$$(h_j^k)^{(t)} \leftarrow \frac{\sum_{i \in \Omega_j} \{r_j^i + (w_k^i)^{(t-1)} (h_j^k)^{(t-1)}\} (w_k^i)^{(t-1)}}{\lambda + \sum_{i \in \Omega_j} \{(w_k^i)^{(t-1)}\}^2}, \quad (3)$$

where  $r_j^i = a_j^i - (\mathbf{w}^i)^{(t-1)}(\mathbf{h}_j)^{(t-1)}$  for all  $(i, j) \in \Omega$ , and a similar rule holds for  $\mathbf{W}$ .

**STRADS implementation:** Our **schedule** strategy is to partition the rows of  $\mathbf{A}$  into  $U$  disjoint index sets  $q_p$ , and the columns of  $\mathbf{A}$  into  $U$  disjoint index sets  $r_p$ . We then dispatch the model variables  $\mathbf{W}, \mathbf{H}$  in round-robin fashion, according to these sets  $q_p, r_p$ . To update elements of  $\mathbf{W}$ , each worker  $p$  computes partial updates on its assigned columns  $r_p$  of  $\mathbf{A}$  and  $\mathbf{H}$ , and analogously for  $\mathbf{H}$  and rows  $q_p$  of  $\mathbf{A}$  and  $\mathbf{W}$ . The sets  $q_p, r_p$  also tie neatly into data partitioning: we merely have to divide  $\mathbf{A}$  into  $U$  pairs of submatrices (where  $U$  is the number of workers), and store the submatrices  $\mathbf{A}^{q_p}$  and  $\mathbf{A}_{r_p}$  at the  $p$ -th worker.

Consider the **push** update for  $\mathbf{H}$  (the case for  $\mathbf{W}$  is similar). To parallel-update a specific element  $(h_j^k)^{(t)}$ , we need  $(w_k^i)^{(t-1)}$  for all  $i \in \Omega_j$ , and  $(\mathbf{h}_j)^{(t-1)}$ . We then compute

$$(a_j^k)_p^{(t)} \leftarrow g_1(k, j, p) := \sum_{i \in (\Omega_j)_p} \left\{ r_j^i + (w_k^i)^{(t-1)}(h_j^k)^{(t-1)} \right\} (w_k^i)^{(t-1)},$$

$$(b_j^k)_p^{(t)} \leftarrow g_2(k, j, p) := \sum_{i \in (\Omega_j)_p} \left\{ (w_k^i)^{(t-1)} \right\}^2,$$

where  $\Omega_j$  are the (observed) elements of column  $\mathbf{A}_j$  in worker  $p$ 's row-submatrix  $\mathbf{A}^{q_p}$ . Finally, **pull** aggregates the updates:

$$(h_j^k)^{(t)} \leftarrow g_3(k, j, [(a_j^k)_p^{(t)}, (b_j^k)_p^{(t)}]) := \frac{\sum_{p=1}^U (a_j^k)_p^{(t)}}{\lambda + \sum_{p=1}^U (b_j^k)_p^{(t)}},$$

with a similar definition for updating  $\mathbf{W}$  using  $(w_k^i)^{(t)} \leftarrow f_3()$  and  $f_1(i, k, p), f_2(i, k, p)$ . This **push-pull** scheme is free from parallelization error: when  $\mathbf{W}$  are updated by **push**, they are mutually independent because  $\mathbf{H}$  is held fixed, and vice-versa. Figure 6 shows the STRADS MF pseudocode.

```
// STRADS Matrix Factorization
```

```
schedule() {  
  // Round-robin scheduling  
  if counter <= U    // Do W  
    return W[q_counter]  
  else                // Do H  
    return H[r_(counter-U)]  
}
```

```
push(worker = p, vars = X[s]) {  
  z = []           // Empty list  
  if counter <= U // X is from W  
    for row in s, k=1..K  
      z.append( (f_1(row,k,p), f_2(row,k,p)) )  
  else            // X is from H  
    for col in s, k=1..K  
      z.append( (g_1(k,col,p), g_2(k,col,p)) )  
  return z  
}
```

```
pull(workers=[p], vars=X[s], updates=[z]) {  
  if counter <= U // X is from W  
    for row in s, k=1..K  
      W[row,k] = f_3(row,k,[z])  
  else            // X is from H  
    for col in s, k=1..K  
      H[k,col] = g_3(k,col,[z])  
  counter = (counter mod 2*U) + 1  
}
```

Figure 6: **STRADS MF pseudocode**. Definitions for  $f_1, g_1, \dots$  and  $q_p, r_p$  are in the text. `counter` is a global model variable.

### 3.3 Lasso

STRADS not only supports simple static **schedules**, but also dynamic, adaptive strategies that take the model state into consideration. Consider Lasso regression [19], which discovers a small

subset of features/dimensions that predict the output  $\mathbf{y}$ . While Lasso can be solved by random parallelization over each dimension’s coefficients, this strategy fails to converge in the presence of strong dependencies between dimensions [4]. Our STRADS Lasso implementation tackles this challenge by (1) avoiding the simultaneous update of coefficients whose dimensions are highly inter-dependent, and (2) prioritizing coefficients that contribute the most to algorithm convergence. These properties complement each other in an algorithmically efficient way, as we shall see.

Formally, Lasso can be defined as an optimization problem:

$$\min_{\boldsymbol{\beta}} \ell(\mathbf{X}, \mathbf{y}, \boldsymbol{\beta}) + \lambda \sum_j |\beta_j|, \quad (4)$$

where  $\lambda$  is a regularization parameter that determines the sparsity of  $\boldsymbol{\beta}$ , and  $\ell(\cdot)$  is a non-negative convex loss function such as squared-loss or logistic-loss; we assume that  $\mathbf{X}$  and  $\mathbf{y}$  are standardized and consider (4) without an intercept. For simplicity but without loss of generality, we let  $\ell(\mathbf{X}, \mathbf{y}, \boldsymbol{\beta}) = \frac{1}{2} \|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|_2^2$ , and note that it is straightforward to use other loss functions. Lasso can be solved using coordinate descent (CD) updates [8]; by taking the gradient of (4), we obtain the CD update rule for  $\beta_j$ :

$$\beta_j^{(t)} \leftarrow S(\mathbf{x}_j^T \mathbf{y} - \sum_{k \neq j} \mathbf{x}_j^T \mathbf{x}_k \beta_k^{(t-1)}, \lambda), \quad (5)$$

where  $S(\cdot, \lambda)$  is a soft-thresholding operator [8], defined by  $S(\beta_j, \lambda) \equiv \text{sign}(\beta) (|\beta| - \lambda)$ .

**STRADS implementation:** Our Lasso **schedule** strategy picks variables dynamically, according to the model state. First, we define a probability distribution  $\mathbf{c} = [c_1, \dots, c_j]$  over the  $\boldsymbol{\beta}$ ; the purpose of  $\mathbf{c}$  is to prioritize  $\beta_j$ ’s during **schedule**, and thus speed up convergence. In particular, we observe that choosing  $\beta_j$  with probability  $c_j = f_1(j) \propto |\beta_j^{(t_j-2)} - \beta_j^{(t_j-1)}| + \eta$  substantially speeds up the Lasso convergence rate (see supplement for our theoretical motivation), where  $\eta$  is a small positive constant, and  $t_j$  is the iteration counter for the  $j$ -th variable.

To prevent non-convergence due to dimension inter-dependencies [4], we only **schedule**  $\beta_j$  and  $\beta_k$  for concurrent updates if  $\mathbf{x}_j^T \mathbf{x}_k \approx 0$ . This is performed as follows: first, select  $U'$  candidates  $\beta_j$ ’s

from the probability distribution  $\mathbf{c}$  to form a set  $\mathcal{C}$ . Next, choose a subset  $\mathcal{B} \subseteq \mathcal{C}$  of size  $U \leq U'$  such that  $\mathbf{x}_j^T \mathbf{x}_k < \rho$  for all  $j, k \in \mathcal{B}$ , where  $\rho \in (0, 1]$ ; we represent this selection procedure<sup>1</sup> by the function  $f_2(\mathcal{C})$ . Here  $U'$  and  $\rho$  are user-defined parameters. We will show that this **schedule** with sufficiently large  $U'$  and small  $\rho$  greatly speeds up convergence over naive random scheduling.

Finally, we execute **push** and **pull** to update the  $\{\beta_j\} \in \mathcal{B}$  using  $U$  workers in parallel. The rows of the data matrix  $\mathbf{X}$  are partitioned into  $U$  submatrices, and the  $p$ -th worker stores the submatrix  $\mathbf{X}^p$ ; With  $\mathbf{X}$  partitioned in this manner, we need to modify the update rule Eq. (5) accordingly. Using  $U$  workers, **push** computes  $U$  partial summations for each selected  $\beta_j$ , denoted by  $\{z_{j,1}^{(t)}, \dots, z_{j,U}^{(t)}\}$ , where  $z_{j,p}^{(t)}$  represents the partial summation for the  $j$ -th  $\beta$  in the  $p$ -th worker at the  $t$ -th iteration:

$$z_{j,p}^{(t)} \leftarrow f_3(p, j) := (\mathbf{x}_j^p)^T \mathbf{y} - \sum_{k \neq j} (\mathbf{x}_j^p)^T (\mathbf{x}_k^p) \beta_k^{(t-1)} \quad (6)$$

After all **pushes** have been completed, **pull** updates  $\beta_j$  via  $\beta_j^{(t)} = f_4(j, [z_{j,p}^{(t)}]) := S(\sum_{p=1}^U z_{j,p}^{(t)}, \lambda)$ .

Figure 7 illustrates the STRADS LASSO pseudocode.

---

<sup>1</sup> Note that this procedure is inexpensive: by selecting  $U'$  candidate  $\beta$ 's first, only  $U'^2$  dependencies need to be checked, as opposed to  $J^2$  where  $J$  is the total number of  $\beta$ .

```
// STRADS Lasso
```

```
schedule() {  
  // Priority-based scheduling  
  for all j      // Get new priorities  
    c_j = f_1(j)  
  for a=1..U'   // Prioritize betas  
    random draw s_a using [c_1, ..., c_j]  
  // Get 'safe' betas  
  (j_1, ..., j_U) = f_2(s_1, ..., s_U')  
  return (b[j_1], ..., b[j_U])  
}
```

```
push(worker = p, vars = (b[j_1],...,b[j_U])) {  
  z = []      // Empty list  
  for a=1..U  // Compute partial sums  
    z.append( f_3(p, j_a) )  
  return z  
}
```

```
pull(workers = [p], vars = (b[j_1],...,b[j_U]),  
      updates = [z]) {  
  for a=1..U      // Aggregate partial sums  
    b[j_a] = f_4(j_a, [z])  
}
```

Figure 7: **STRADS Lasso pseudocode.** Definitions for  $f_1, f_2, \dots$  are given in the text.

## 4. EXPERIMENTS

We now demonstrate that our STRADS implementations of LDA, MF and Lasso can (1) reach larger model sizes than other baselines; (2) converge at least as fast, if not faster, than other baselines; (3) with additional machines, STRADS uses less memory per machine (efficient partitioning). For baselines, we used (a) a STRADS implementation of distributed Lasso with only a naive round-robin scheduler (Lasso-RR), (b) GraphLab’s Alternating Least Squares (ALS) implementation of MF [14], (c) YahooLDA for topic modeling [1]. Note that Lasso-RR imitates the



random scheduling scheme proposed by Shotgun algorithm on STRADS. We chose GraphLab and YahooLDA, as they are popular choices for distributed MF and LDA.

We conducted experiments on two clusters [10] (with 2-core and 16-core machines respectively), to show the effectiveness of STRADS model-parallelism across different hardware. We used the 2-core cluster for LDA, and the 16-core cluster for Lasso and MF. The 2-core cluster contains 128 machines, each with two 2.6GHz AMD cores and 8GB RAM, and connected via a 1Gbps network interface. The 16-core cluster contains 9 machines, each with 16 2.1GHz AMD cores and 64GB RAM, and connected via a 40Gbps network interface. All our experiments use a fixed data size, and we vary the number of machines and/or the model size (unless otherwise stated).

#### 4.1 Datasets

**Latent Dirichlet Allocation** We used 3.9M English Wikipedia abstracts, and conducted experiments using both unigram (1-word) tokens ( $V = 2.5\text{M}$  unique unigrams, 179M tokens) and bigram (2-word) tokens ( $V = 21.8\text{M}$  unique bigrams, 79M tokens). We note that our bigram vocabulary (21.8M) is an order of magnitude larger than recently published results [1], demonstrating that STRADS scales to very large models. We set the number of topics to  $K = 5000$  and 10000 (again, significantly larger than recent literature [1]), which creates extremely large word-topic tables: 12.5B elements (unigram) and 109B elements (bigram).

**Matrix Factorization** We used the Netflix dataset [2] for our MF experiments: 100M anonymized ratings from 480,189 users on 17,770 movies. We varied the rank of  $\mathbf{W}$ ,  $\mathbf{H}$  from  $K = 20$  to 2000, which exceeds the upper limit of previous MF papers [23, 9, 21].

**Lasso** We used synthetic data with 50K samples and  $J = 10\text{M}$  to 100M features, where every feature  $\mathbf{x}_j$  has only 25 non-zero samples. To simulate correlations between adjacent features (which exist in real-world data), we first added  $Unif(0, 1)$  noise to  $\mathbf{x}_1$ . Then, for  $j = 2, \dots, J$ , with 0.9 probability we add  $\epsilon_j = Unif(0, 1)$  noise to  $\mathbf{x}_j$ , otherwise we add  $0.9\epsilon_{j-1} + 0.1Unif(0, 1)$  to  $\mathbf{x}_j$ .

## 4.2 Speed and Model Sizes

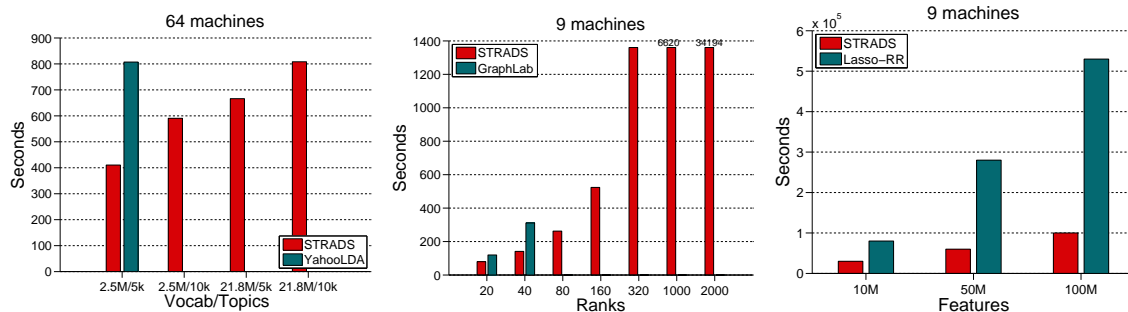


Figure 8: Convergence time versus model size for STRADS and baselines for (left) LDA, (center) MF, and (right) Lasso. We omit the bars if a method did not reach 98% of STRADS’s convergence point (YahooLDA and GraphLab-MF failed at 2.5M-Vocab/10K-topics and rank  $K \geq 80$ , respectively). STRADS not only reaches larger model sizes than YahooLDA, GraphLab, and Lasso-RR, but also converges significantly faster.

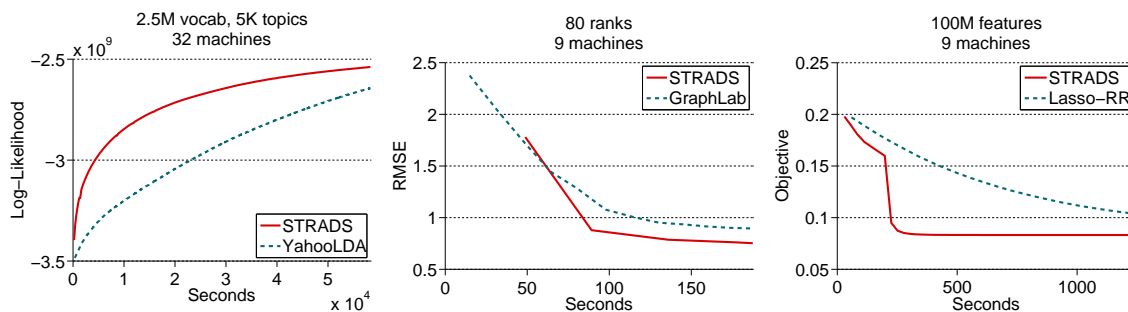


Figure 9: Convergence trajectories of different methods for (left) LDA, (center) MF, and (right) Lasso.

Figure 8 shows the time taken by each algorithm to reach a fixed objective value (over a range of model sizes), as well as the largest model size that each baseline was capable of running. For LDA and MF, STRADS handles much larger model sizes than either YahooLDA (could only handle 5K topics on the unigram dataset) or GraphLab (could only handle rank  $< 80$ ), while converging more quickly; we attribute STRADS’s faster convergence to lower parallelization error (LDA only) and reduced synchronization requirements through careful model partitioning (LDA, MF). In particular, YahooLDA stores nearly the whole word-topic table on every machine, so its maximum model size is limited by the smallest machine (Figure 3). For Lasso, STRADS converges

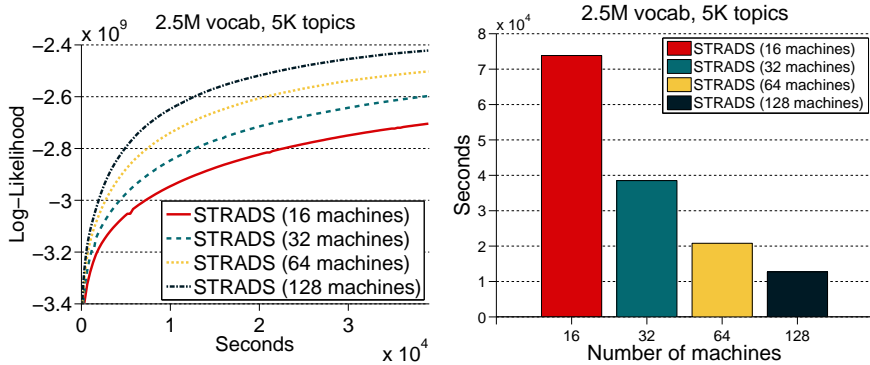


Figure 10: STRADS LDA scalability with increasing machines using a fixed model size. (Left) Convergence trajectories; (Right) Time taken to reach a log-likelihood of  $-2.6 \times 10^9$ .

more quickly than Lasso-RR because of our dynamic **schedule** strategy, which is graphically captured in the convergence trajectory seen in Figure 9 — observe that STRADS’s dynamic **schedule** causes the Lasso objective to plunge quickly to the optimum at around 250 seconds. We also see that STRADS LDA and MF achieved better objective values, confirming that STRADS model-parallelism is fast without compromising convergence quality.

### 4.3 Scalability

In Figure 10, we show the convergence trajectories and time-to-convergence for STRADS LDA using different numbers of machines at a fixed model size (unigram with 2.5M vocab and 5K topics). The plots confirm that STRADS LDA exhibits faster convergence with more machines, and that the time to convergence almost halves with every doubling of machines (near-linear scaling).

## 5. DISCUSSION AND RELATED WORK

As a framework of user-programmable primitives for dynamic Big Model-parallelism, STRADS provides the following benefits: (1) scalability and efficient memory utilization, allowing larger models to be run with additional machines (because the model is partitioned, rather than duplicated across machines); (2) the ability to invoke dynamic **schedules** that reduce model variable dependencies across workers, leading to lower parallelization error and thus faster, correct convergence.

While the notion of model-parallelism is not new, our contribution is to study it within the context of a programmable system (STRADS), using primitives that enable general, user-programmable partitioning and static/dynamic scheduling of variable updates (based on model dependencies). Previous works explore aspects of model-parallelism in a more specific context: Scherrer et al. [18] proposed a static model partitioning scheme specifically for parallel coordinate descent, while GraphLab [15, 14] statically pre-partitions data and variables through a graph abstraction.

An important direction for future research is to reduce the communication costs of using STRADS. Currently, STRADS adopts a star topology from scheduler machines to workers, which causes the scheduler to eventually become a bottleneck as we increase the number of machines. To mitigate this issue, we wish to explore different **sync** schemes such as an asynchronous parallelism [1] and stale synchronous parallelism [13]. We also want to explore the use of STRADS for other popular ML applications, such as support vector machines and logistic regression.

## REFERENCES

- [1] Amr Ahmed, Moahmed Aly, Joseph Gonzalez, Shravan Narayanamurthy, and Alexander J Smola. Scalable inference in latent variable models. In *WSDM*, pages 123–132. ACM, 2012.
- [2] James Bennett and Stan Lanning. The netflix prize. In *Proceedings of KDD cup and workshop*, volume 2007, page 35, 2007.
- [3] David M Blei, Andrew Y Ng, and Michael I Jordan. Latent dirichlet allocation. *the Journal of machine Learning research*, 3:993–1022, 2003.
- [4] Joseph K Bradley, Aapo Kyrola, Danny Bickson, and Carlos Guestrin. Parallel coordinate descent for l1-regularized loss minimization. *ICML*, 2011.
- [5] J. Dean, G. Corrado, R. Monga, K. Chen, M. Devin, Q. V. Le, M. Z. Mao, M. Ranzato, A. W. Senior, P. A. Tucker, et al. Large scale distributed deep networks. In *NIPS*, pages 1232–1240, 2012.
- [6] Jeffrey Dean and Sanjay Ghemawat. MapReduce: simplified data processing on large clusters. *Communications of the ACM*, 51(1):107–113, 2008.

- [7] Jianqing Fan, Richard Samworth, and Yichao Wu. Ultrahigh dimensional feature selection: beyond the linear model. *The Journal of Machine Learning Research*, 10:2013–2038, 2009.
- [8] J. Friedman, T. Hastie, H. Hofling, and R. Tibshirani. Pathwise coordinate optimization. *Annals of Applied Statistics*, 1(2):302–332, 2007.
- [9] Rainer Gemulla, Erik Nijkamp, Peter J Haas, and Yannis Sismanis. Large-scale matrix factorization with distributed stochastic gradient descent. In *Proceedings of the 17th ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 69–77. ACM, 2011.
- [10] Garth Gibson, Gary Grider, Andree Jacobson, and Wyatt Lloyd. Probe: A thousand-node experimental cluster for computer systems research. *USENIX; login*, 38, 2013.
- [11] J. Gonzalez, Y. Low, A. Gretton, and C. Guestrin. Parallel gibbs sampling: From colored fields to thin junction trees. In *International Conference on Artificial Intelligence and Statistics*, pages 324–332, 2011.
- [12] Thomas L Griffiths and Mark Steyvers. Finding scientific topics. *Proceedings of the National Academy of Sciences of the United States of America*, 101(Suppl 1):5228–5235, 2004.
- [13] Q. Ho, J. Cipar, H. Cui, J.-K. Kim, S. Lee, P. B. Gibbons, G. Gibson, G. R. Ganger, and E. P. Xing. More effective distributed ml via a stale synchronous parallel parameter server. In *NIPS*, 2013.
- [14] Y. Low, J. Gonzalez, A. Kyrola, D. Bickson, C. Guestrin, and J. M. Hellerstein. Distributed GraphLab: A Framework for Machine Learning and Data Mining in the Cloud. *PVLDB*, 2012.
- [15] Yucheng Low, Joseph Gonzalez, Aapo Kyrola, Danny Bickson, Carlos Guestrin, and Joseph M. Hellerstein. Graphlab: A new parallel framework for machine learning. In *UAI*, July 2010.
- [16] Grzegorz Malewicz, Matthew H Austern, Aart JC Bik, James C Dehnert, Ilan Horn, Naty Leiser, and Grzegorz Czajkowski. Pregel: a system for large-scale graph processing. In *Proceedings of the 2010 ACM SIGMOD International Conference on Management of data*, pages 135–146. ACM, 2010.
- [17] David Newman, Arthur Asuncion, Padhraic Smyth, and Max Welling. Distributed algorithms for topic models. *The Journal of Machine Learning Research*, 10:1801–1828, 2009.

- [18] Chad Scherrer, Ambuj Tewari, Mahantesh Halappanavar, and David Haglin. Feature clustering for accelerating parallel coordinate descent. *NIPS*, 2012.
- [19] R. Tibshirani. Regression shrinkage and selection via the lasso. *Journal of the Royal Statistical Society. Series B (Methodological)*, 58(1):267–288, 1996.
- [20] Limin Yao, David Mimno, and Andrew McCallum. Efficient methods for topic model inference on streaming document collections. 2009.
- [21] Hsiang-Fu Yu, Cho-Jui Hsieh, Si Si, and Inderjit Dhillon. Scalable coordinate descent approaches to parallel matrix factorization for recommender systems. In *ICDM 2012*, pages 765–774. IEEE, 2012.
- [22] M. Zaharia, M. Chowdhury, M. J. Franklin, S. Shenker, and I. Stoica. Spark: cluster computing with working sets. In *Proceedings of the 2nd USENIX conference on Hot topics in cloud computing*, 2010.
- [23] Y. Zhou, D. Wilkinson, R. Schreiber, and R. Pan. Large-scale parallel collaborative filtering for the netflix prize. In *Algorithmic Aspects in Information and Management*, pages 337–348. Springer, 2008.
- [24] Martin Zinkevich, John Langford, and Alex J Smola. Slow learners are fast. In *Advances in Neural Information Processing Systems*, pages 2331–2339, 2009.
- [25] Martin Zinkevich, Markus Weimer, Lihong Li, and Alex J Smola. Parallelized stochastic gradient descent. In *Advances in Neural Information Processing Systems*, pages 2595–2603, 2010.