

# PASSCoDe: Parallel ASynchronous Stochastic dual Co-ordinate Descent

Cho-Jui Hsieh

Department of Computer Science  
University of Texas at Austin

Joint work with H.-F. Yu and I. S. Dhillon

# Outline

- L2-regularized Empirical Risk Minimization
- Dual Coordinate Descent (Hsieh et al., 2008)
- Parallel Dual Coordinate Descent (on multi-core machines)
- Theoretical Analysis
- Experimental Results

# L2-regularized ERM

$$\mathbf{w}^* = \arg \min_{\mathbf{w} \in R^d} P(\mathbf{w}) := \frac{1}{2} \|\mathbf{w}\|^2 + \sum_{i=1}^n \ell_i(\mathbf{w}^T \mathbf{x}_i)$$

- SVM with hinge loss:  $\ell_i(z_i) = C \max(1 - z_i, 0)$
- SVM with squared hinge loss:  $\ell_i(z_i) = C \max(1 - z_i, 0)^2$
- Logistic regression:  $\ell_i(z_i) = C \log(1 + e^{-z_i})$

# Primal and Dual Formulations

- Primal Problem

$$\mathbf{w}^* = \arg \min_{\mathbf{w} \in \mathbb{R}^d} P(\mathbf{w}) := \frac{1}{2} \|\mathbf{w}\|^2 + \sum_{i=1}^n \ell_i(\mathbf{w}^T \mathbf{x}_i)$$

- Dual Problem

$$\alpha^* = \arg \min_{\alpha \in \mathbb{R}^n} D(\alpha) := \frac{1}{2} \left\| \sum_{i=1}^n \alpha_i \mathbf{x}_i \right\|^2 + \sum_{i=1}^n \ell_i^*(-\alpha_i),$$

- $\ell_i^*(\cdot)$ : the conjugate of  $\ell_i(\cdot)$
- Primal-Dual Relationship between  $\mathbf{w}^*$  and  $\alpha^*$

$$\mathbf{w}^* = \mathbf{w}(\alpha^*) := \sum_{i=1}^n \alpha_i^* \mathbf{x}_i$$

# Coordinate Descent on the Dual Problem

Randomly select an  $i \in \{1, \dots, n\}$  and update  $\alpha_i \leftarrow \alpha_i + \delta^*$ , where

$$\delta^* = \arg \min_{\delta} D(\boldsymbol{\alpha} + \delta \mathbf{e}_i)$$

# Coordinate Descent on the Dual Problem

Randomly select an  $i \in \{1, \dots, n\}$  and update  $\alpha_i \leftarrow \alpha_i + \delta^*$ , where

$$\begin{aligned}\delta^* &= \arg \min_{\delta} D(\boldsymbol{\alpha} + \delta \mathbf{e}_i) \\ &= \arg \min_{\delta} \frac{1}{2} \left( \delta + \frac{(\sum_{i=1}^n \alpha_i \mathbf{x}_i)^T \mathbf{x}_i}{\|\mathbf{x}_i\|^2} \right)^2 + \frac{1}{\|\mathbf{x}_i\|^2} \ell_i^* (-(\alpha_i + \delta)) \\ &= T_i \left( \left( \sum_{i=1}^n \alpha_i \mathbf{x}_i \right)^T \mathbf{x}_i, \alpha_i \right)\end{aligned}$$

- Simple univariate problem, but  $O(nnz)$  construction time

# Coordinate Descent on the Dual Problem

Randomly select an  $i \in \{1, \dots, n\}$  and update  $\alpha_i \leftarrow \alpha_i + \delta^*$ , where

$$\begin{aligned}\delta^* &= \arg \min_{\delta} D(\boldsymbol{\alpha} + \delta \mathbf{e}_i) \\ &= \arg \min_{\delta} \frac{1}{2} \left( \delta + \frac{(\sum_{i=1}^n \alpha_i \mathbf{x}_i)^T \mathbf{x}_i}{\|\mathbf{x}_i\|^2} \right)^2 + \frac{1}{\|\mathbf{x}_i\|^2} \ell_i^* (-(\alpha_i + \delta)) \\ &= T_i \left( \left( \sum_{i=1}^n \alpha_i \mathbf{x}_i \right)^T \mathbf{x}_i, \alpha_i \right)\end{aligned}$$

- Simple univariate problem, but  ~~$O(nnz)$~~  construction time  $\Rightarrow O(n_i)$

## DCD: [Hsieh et al 2008]

- Maintain primal variable  $\mathbf{w} = \sum_{i=1}^n \alpha_i \mathbf{x}_i$  and  $\delta^* = T_i(\mathbf{w}^T \mathbf{x}_i, \alpha_i)$
- $O(n_i)$  construction time:  $n_i = nnz$  of  $\mathbf{x}_i$
- $O(n_i)$  maintenance cost:  $\mathbf{w} \leftarrow \mathbf{w} + \delta^* \mathbf{x}_i$

## Stochastic Dual Coordinate Descent

For  $t = 1, 2, \dots$

- 1 Randomly pick an index  $i$
- 2 Compute  $\mathbf{w}^T \mathbf{x}_i$
- 3 Update  $\alpha_i \leftarrow \alpha_i + \delta^*$  where  $\delta^* = T_i(\mathbf{w}^T \mathbf{x}_i, \alpha_i)$
- 4 Update  $\mathbf{w} \leftarrow \mathbf{w} + \delta^* \mathbf{x}_i$ .

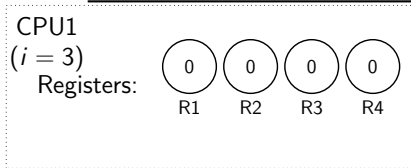
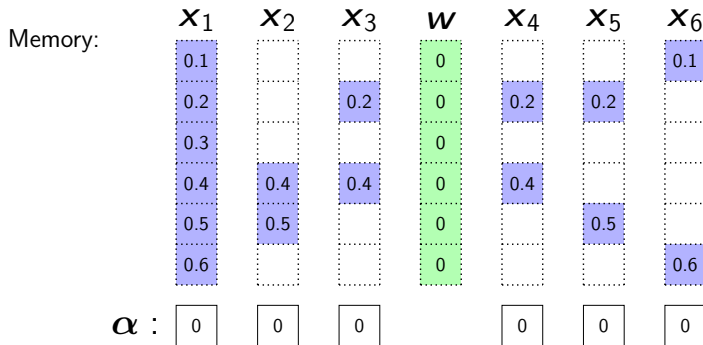
- Implemented in LIBLINEAR:

Linear SVM (Hsieh et al., 2008), multi-class SVM (Keerthi et al., 2008), Logistic regression (Yu et al., 2011).

- Analysis: (Nesterov et al., 2012; Shalev-Shwartz et al., 2013)

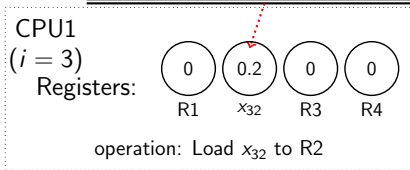
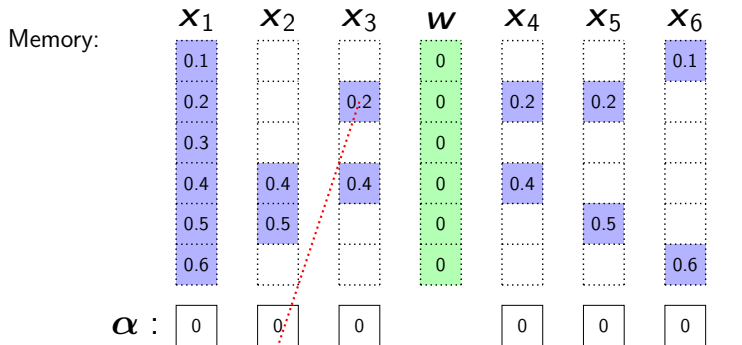


# Dual Coordinate Descent



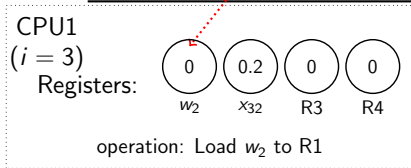
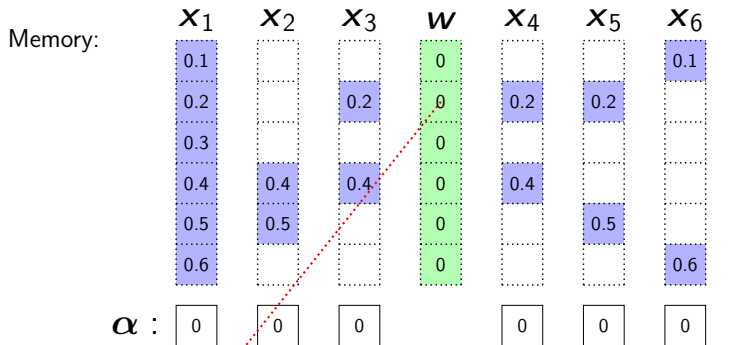
DCD step: compute  $w^T x_i$

# Dual Coordinate Descent



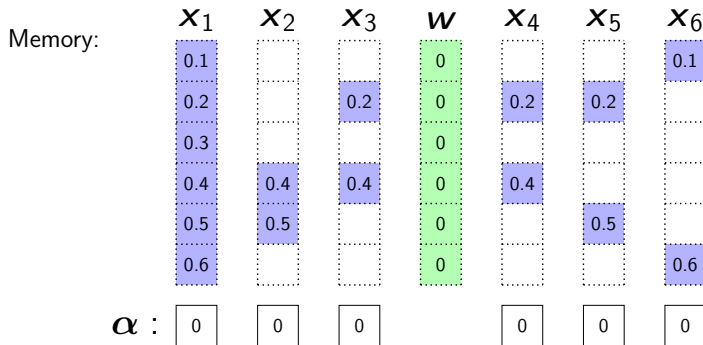
DCD step: compute  $w^T x_i$

# Dual Coordinate Descent



DCD step: compute  $w^T x_i$

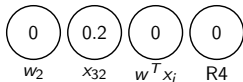
# Dual Coordinate Descent



CPU1

( $i = 3$ )

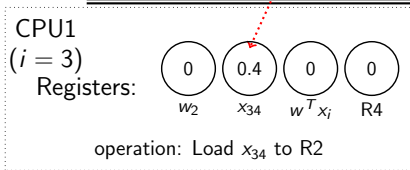
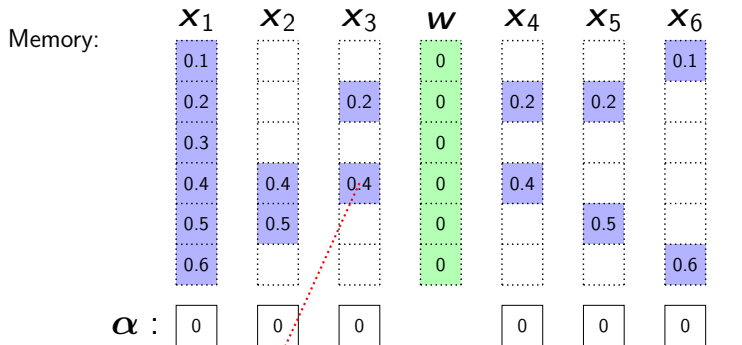
Registers:



operation:  $R_3 = R_3 + R_1 \times R_2$

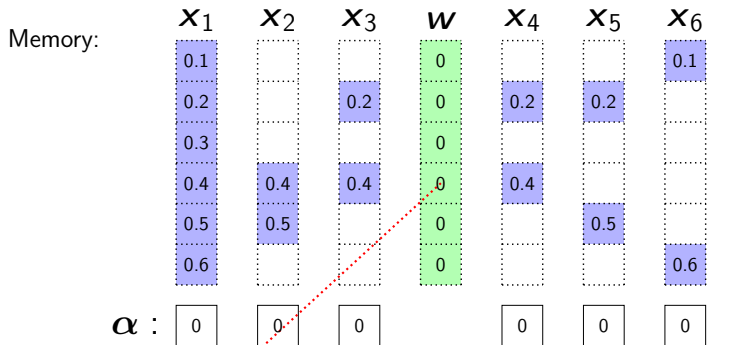
DCD step: compute  $w^T x_i$

# Dual Coordinate Descent



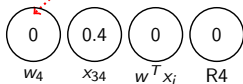
DCD step: compute  $w^T x_i$

# Dual Coordinate Descent



CPU1  
( $i = 3$ )

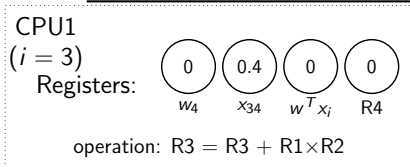
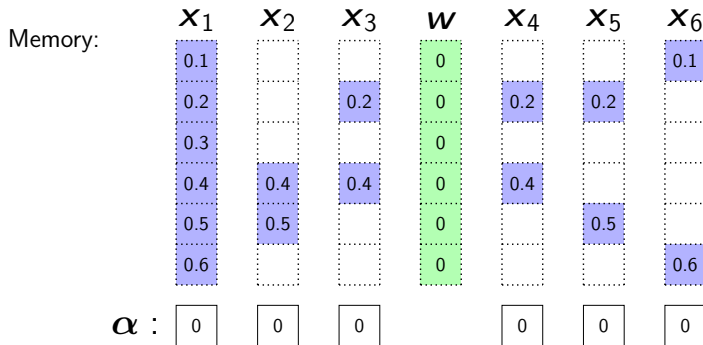
Registers:



operation: Load  $w_4$  to R1

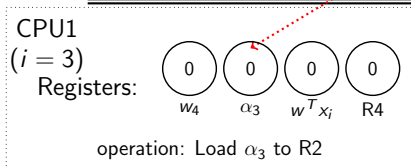
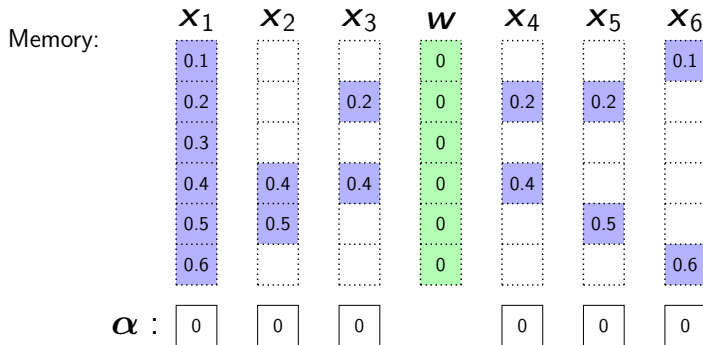
DCD step: compute  $w^T x_i$

# Dual Coordinate Descent



DCD step: compute  $w^T x_i$

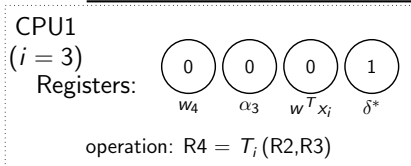
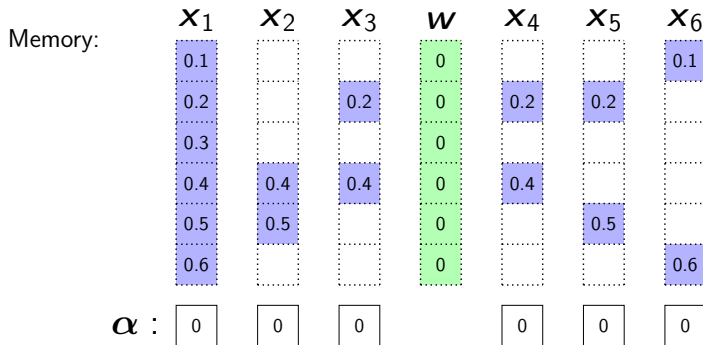
# Dual Coordinate Descent



DCD step: compute  $\delta^* = T_j(w^T x, \alpha_j)$

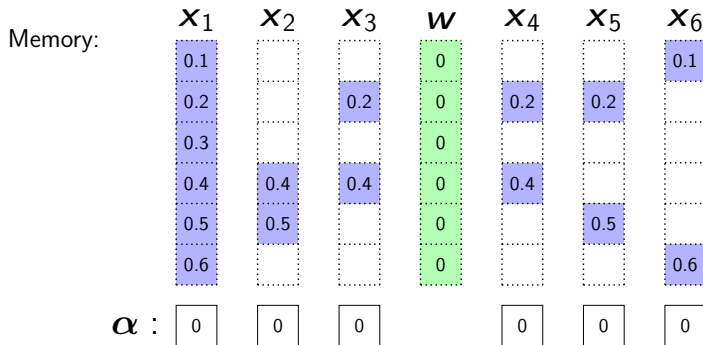


# Dual Coordinate Descent



DCD step: compute  $\delta^* = T_i(w^T x, \alpha_i)$

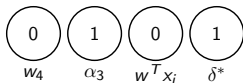
# Dual Coordinate Descent



CPU1

( $i = 3$ )

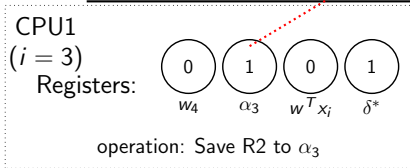
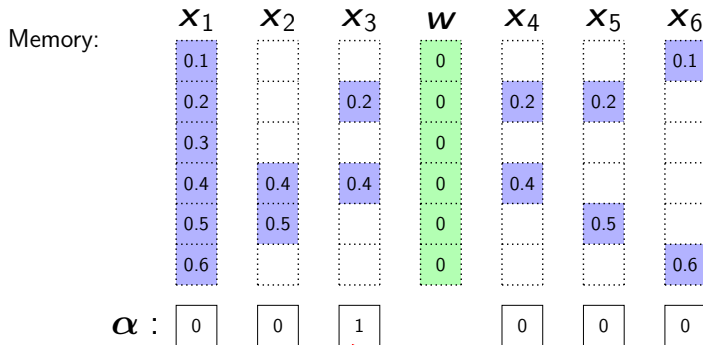
Registers:



operation:  $R2 = R2 + R4$

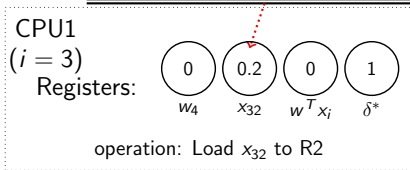
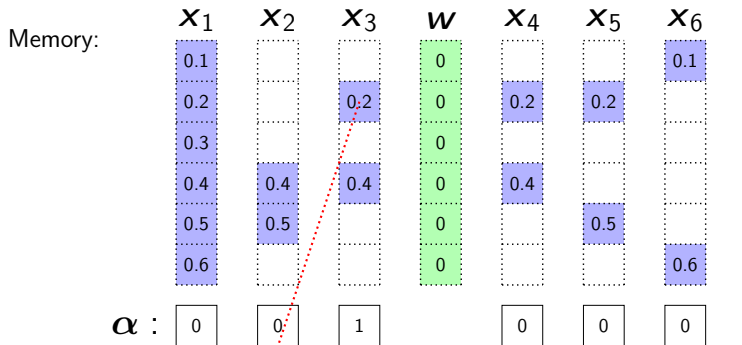
DCD step: **update**  $\alpha_j = \alpha_j + \delta^*$

# Dual Coordinate Descent



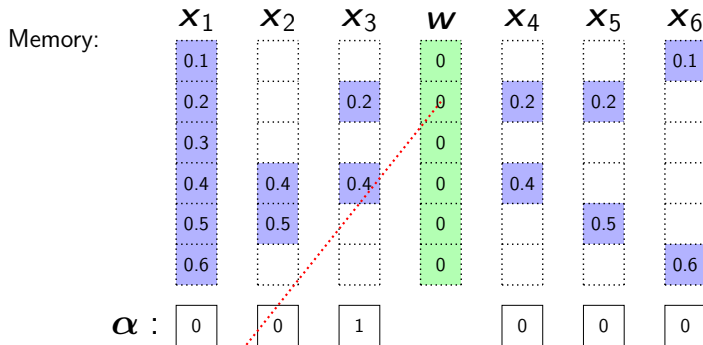
DCD step: update  $\alpha_j = \alpha_j + \delta^*$

# Dual Coordinate Descent



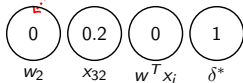
DCD step: update  $w = w + \delta^* x_i$

# Dual Coordinate Descent



CPU1  
( $i = 3$ )

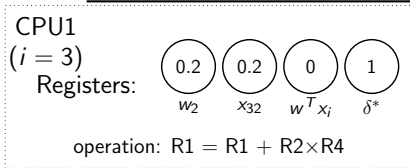
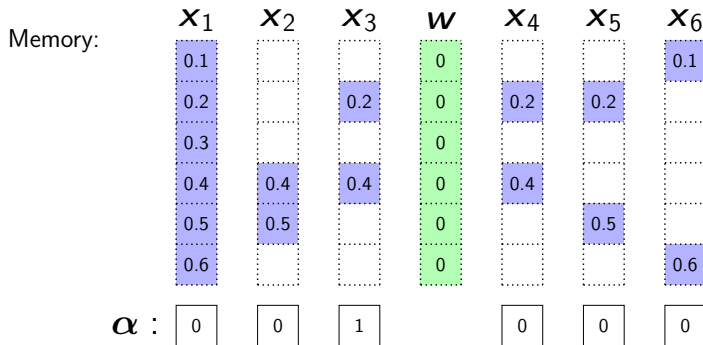
Registers:



operation: Load  $w_2$  to R1

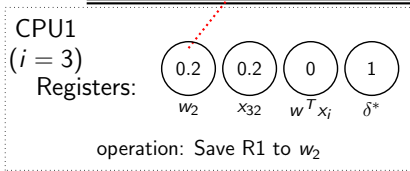
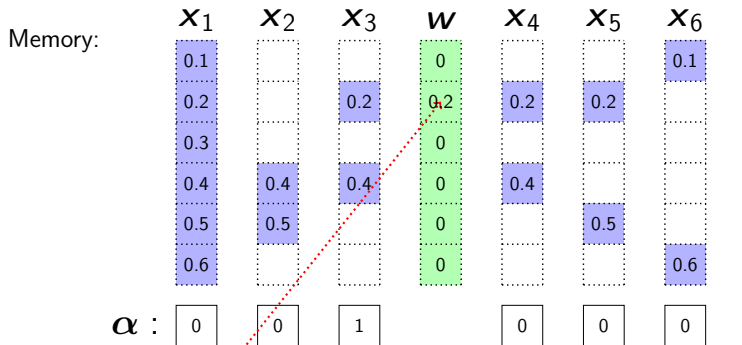
DCD step: update  $w = w + \delta^* x_i$

# Dual Coordinate Descent



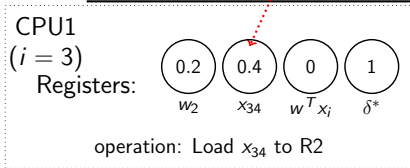
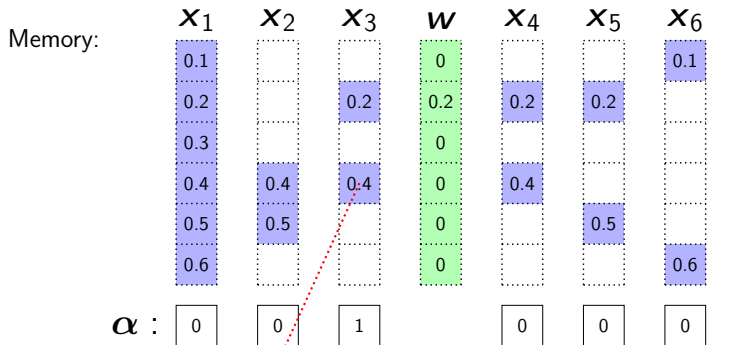
DCD step: update  $w = w + \delta^* x_i$

# Dual Coordinate Descent



DCD step: update  $w = w + \delta^* x_i$

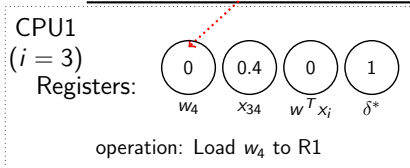
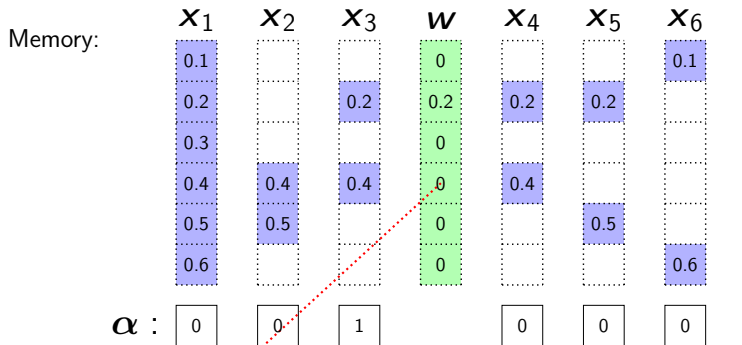
# Dual Coordinate Descent



DCD step: update  $w = w + \delta^* x_i$

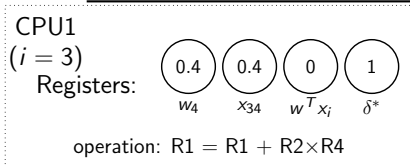
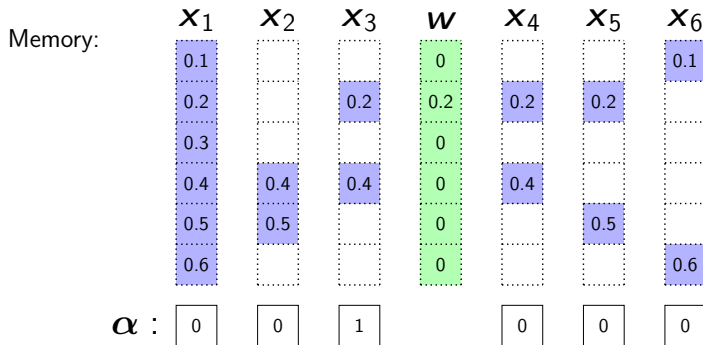


# Dual Coordinate Descent



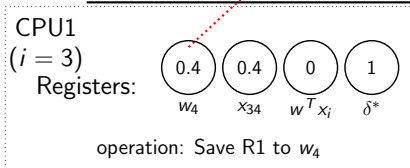
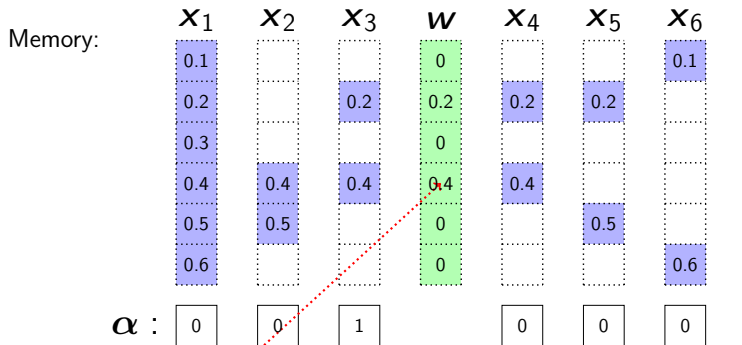
DCD step: update  $w = w + \delta^* x_j$

# Dual Coordinate Descent



DCD step: update  $w = w + \delta^* x_i$

# Dual Coordinate Descent



DCD step: update  $w = w + \delta^* x_i$

# Parallel DCD in Shared-memory Multi-core System

- Serial DCD updates:

For  $t = 1, 2, \dots$

- 1 Randomly pick an index  $i$
- 2 Compute  $\mathbf{w}^T \mathbf{x}_i$
- 3 Update  $\alpha_i \leftarrow \alpha_i + \delta^*$  where  $\delta^* = T_i(\mathbf{w}^T \mathbf{x}_i, \alpha_i)$
- 4 Update  $\mathbf{w} \leftarrow \mathbf{w} + \delta^* \mathbf{x}_i$ .

# Parallel DCD in Shared-memory Multi-core System

- Parallel DCD updates:

**Each thread repeatedly performs the following updates.**

For  $t = 1, 2, \dots$

- 1 Randomly pick an index  $i$
- 2 Compute  $\mathbf{w}^T \mathbf{x}_i$
- 3 Update  $\alpha_i \leftarrow \alpha_i + \delta^*$  where  $\delta^* = T_i(\mathbf{w}^T \mathbf{x}_i, \alpha_i)$
- 4 Update  $\mathbf{w} \leftarrow \mathbf{w} + \delta^* \mathbf{x}_i$ .

# Parallel DCD in Shared-memory Multi-core System

- Parallel DCD updates:

**Each thread repeatedly performs the following updates.**

For  $t = 1, 2, \dots$

- 1 Randomly pick an index  $i$
- 2 Compute  $\mathbf{w}^T \mathbf{x}_i$
- 3 Update  $\alpha_i \leftarrow \alpha_i + \delta^*$  where  $\delta^* = T_i(\mathbf{w}^T \mathbf{x}_i, \alpha_i)$
- 4 Update  $\mathbf{w} \leftarrow \mathbf{w} + \delta^* \mathbf{x}_i$ .

- Easy to implement using OpenMP.
- Variables  $\alpha$  and  $\mathbf{w}$  stored in shared memory.

# Parallel DCD in Shared-memory Multi-core System

- Parallel DCD updates:

**Each thread repeatedly performs the following updates.**

For  $t = 1, 2, \dots$

- 1 Randomly pick an index  $i$
- 2 Compute  $\mathbf{w}^T \mathbf{x}_i$
- 3 Update  $\alpha_i \leftarrow \alpha_i + \delta^*$  where  $\delta^* = T_i(\mathbf{w}^T \mathbf{x}_i, \alpha_i)$
- 4 Update  $\mathbf{w} \leftarrow \mathbf{w} + \delta^* \mathbf{x}_i$ .

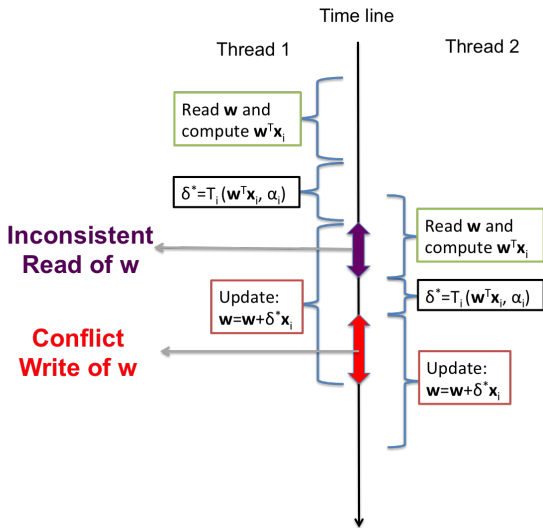
- Easy to implement using OpenMP.
- Variables  $\alpha$  and  $\mathbf{w}$  stored in shared memory.
- Distributed Dual Coordinate Descent:

Each machine has local copy of  $\alpha$ ,  $\mathbf{w}$

(Yang, 2013; Jaggi et al, 2014; Lee and Roth, 2015; Ma et al., 2015).

# Parallel Dual Coordinate Descent: Two Issues for Correctness

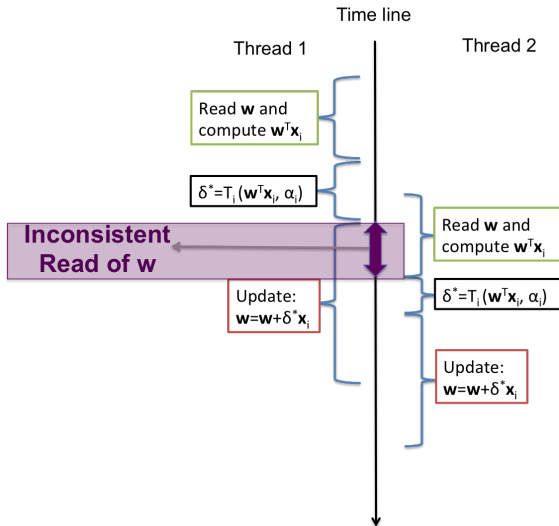
- Inconsistent Read of  $\mathbf{w}$
- Conflict Write of  $\mathbf{w}$





# Inconsistent Read

- Thread 2 reads  $\mathbf{w}$  and thread 1 writes to  $\mathbf{w}$  simultaneously.
- There may not exist any  $\alpha$  such that  $\mathbf{w} = \sum_i \alpha_i \mathbf{x}_i$ .
- The “bounded delay” analysis in (Liu and Wright, 2014) cannot be directly applied.

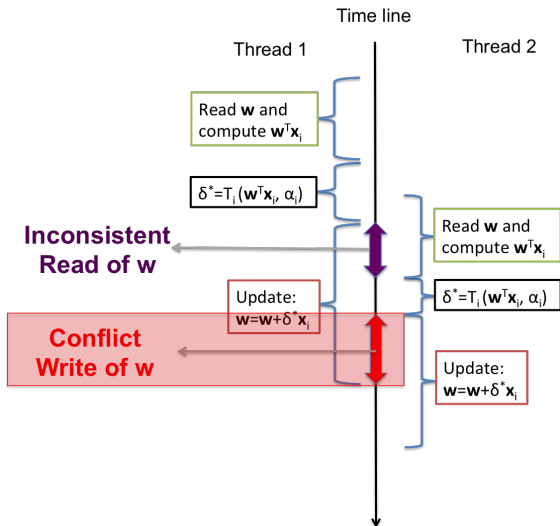


# Conflict Write

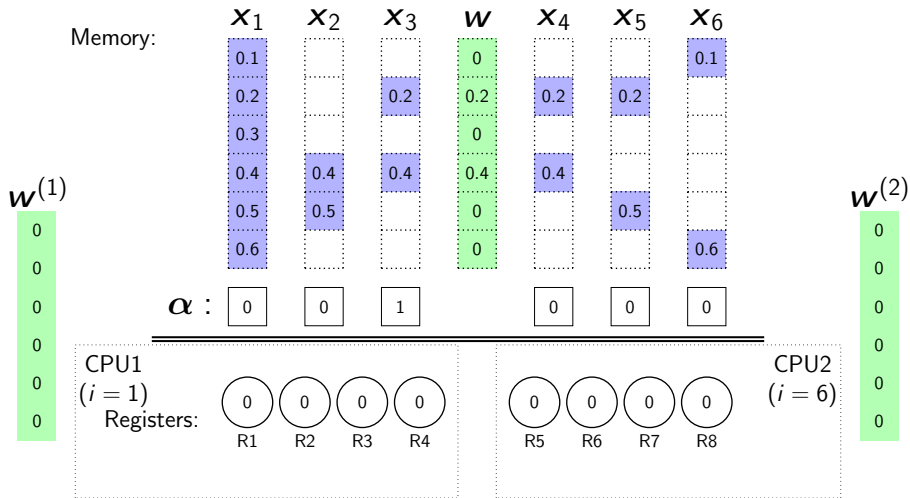
- Thread 1 and 2 write to  $w$  simultaneously.
- Updates to  $w$  can be overwritten, so the converged solution  $\hat{w}$  and  $\hat{\alpha}$  may be inconsistent:

$$\hat{w} \neq \sum_i \hat{\alpha}_i x_i.$$

	CPU1:			CPU2:	
	OP	R1	w	OP	R2
0		0.0	1.0		0.0
1	load w	1.0	1.0	load w	1.0
2	add 0.2	1.2	1.0	add 0.5	1.5
3	save w	1.2	1.2		1.5
4		1.2	1.5	save w	1.5

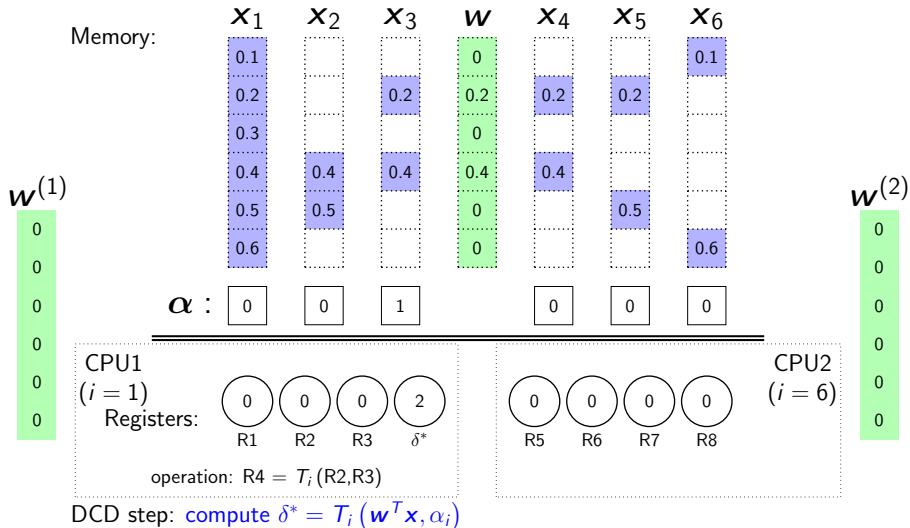


# Dual Coordinate Descent in Parallel

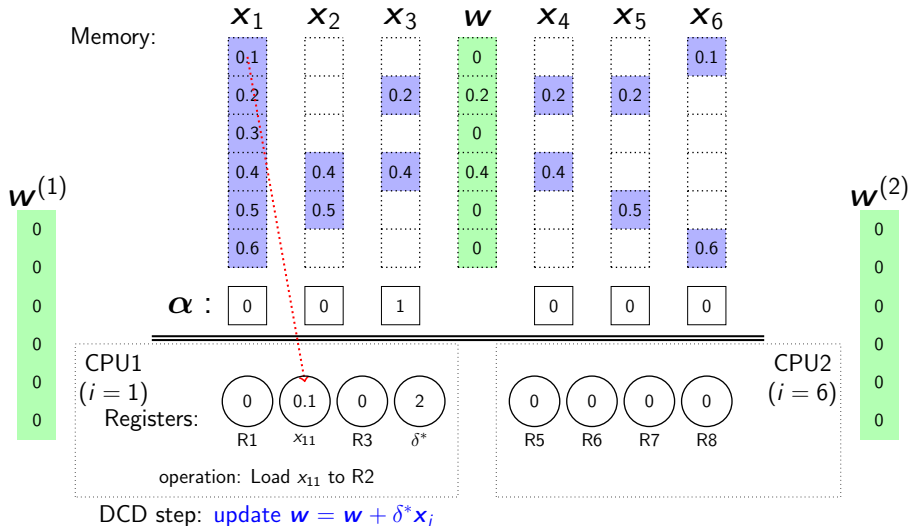


DCD step: compute  $\delta^* = T_i(w^T x, \alpha_i)$

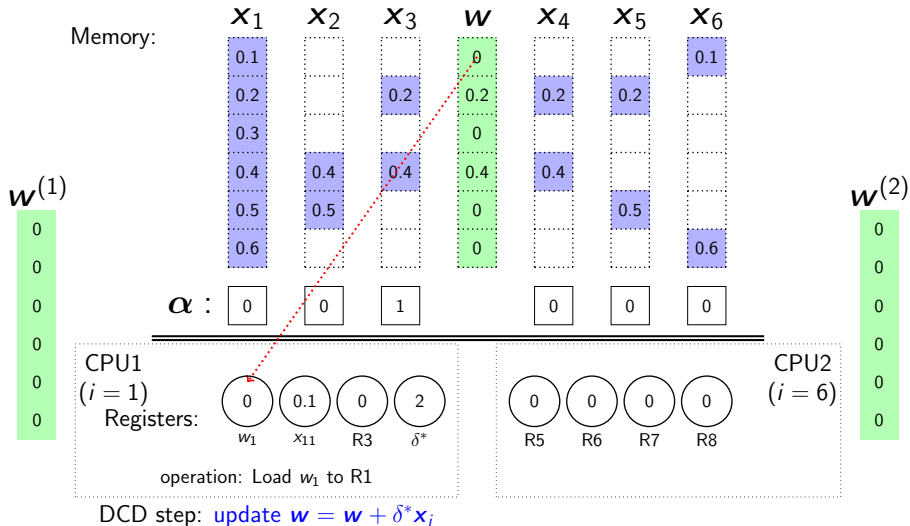
# Dual Coordinate Descent in Parallel



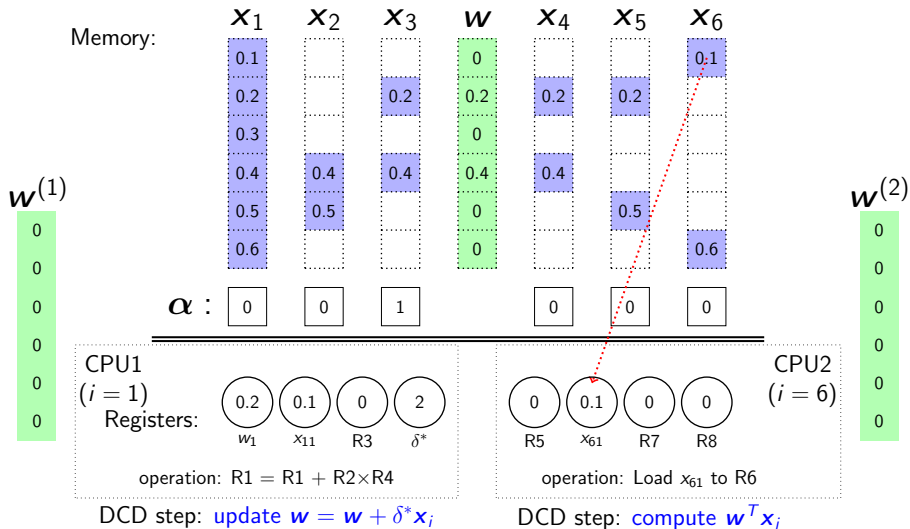
# Dual Coordinate Descent in Parallel



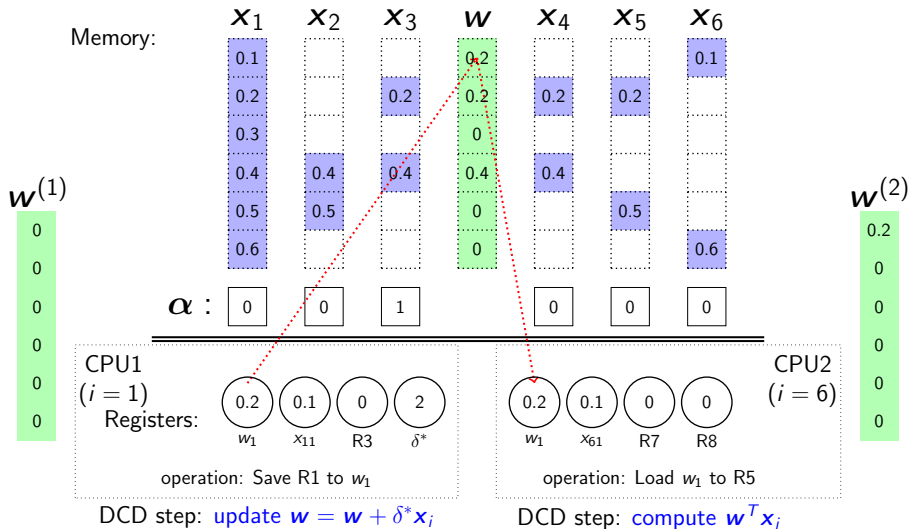
# Dual Coordinate Descent in Parallel



# Dual Coordinate Descent in Parallel

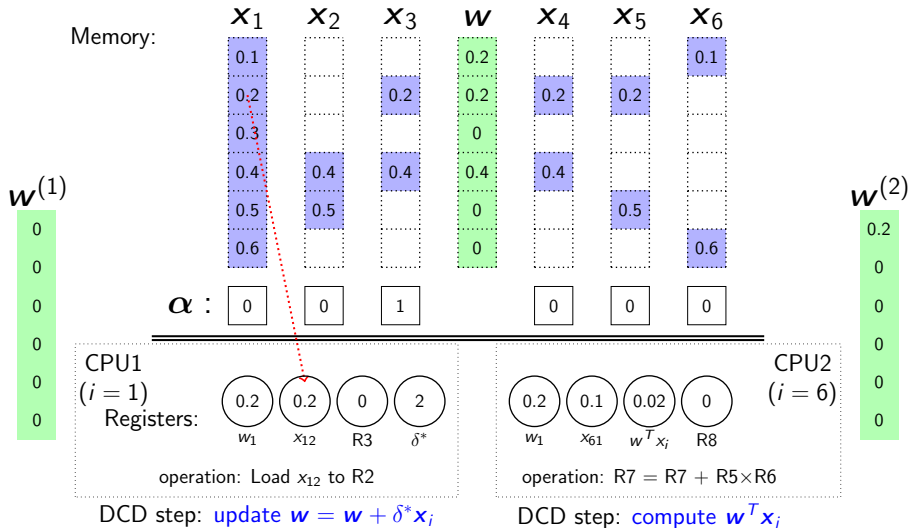


# Dual Coordinate Descent in Parallel

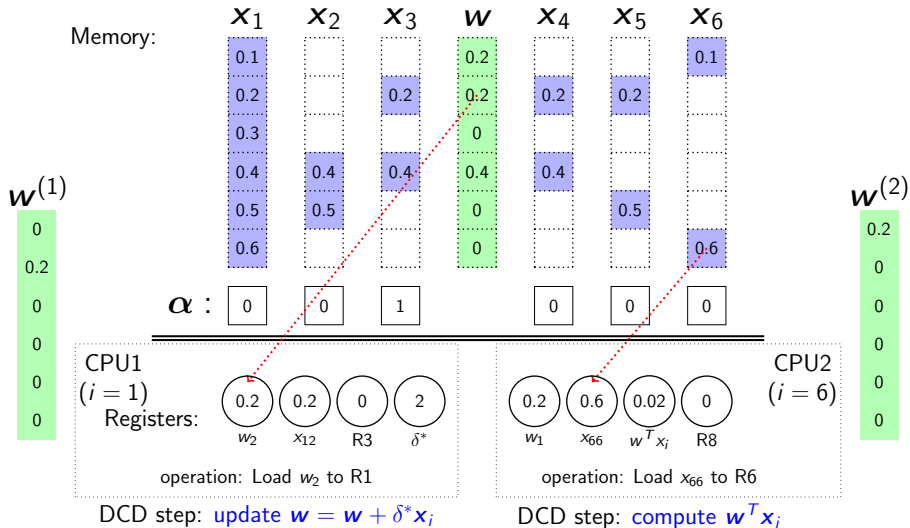




# Dual Coordinate Descent in Parallel

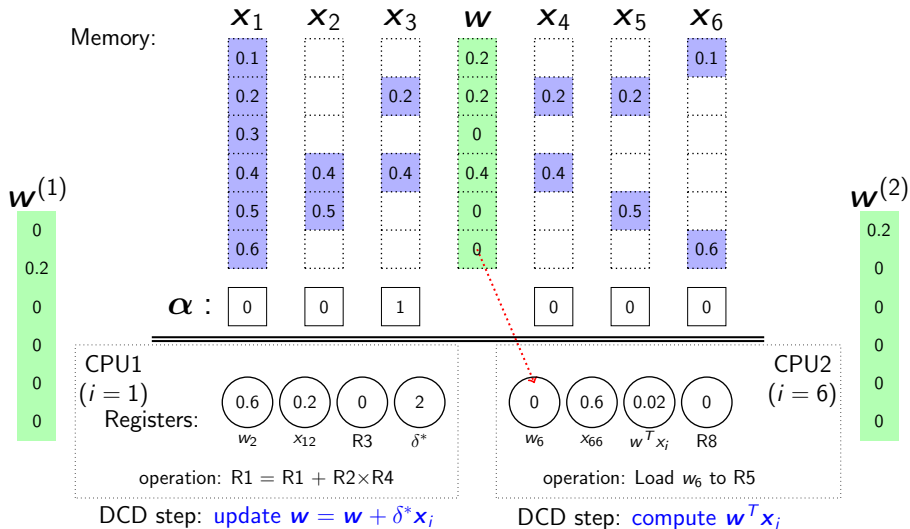


# Dual Coordinate Descent in Parallel



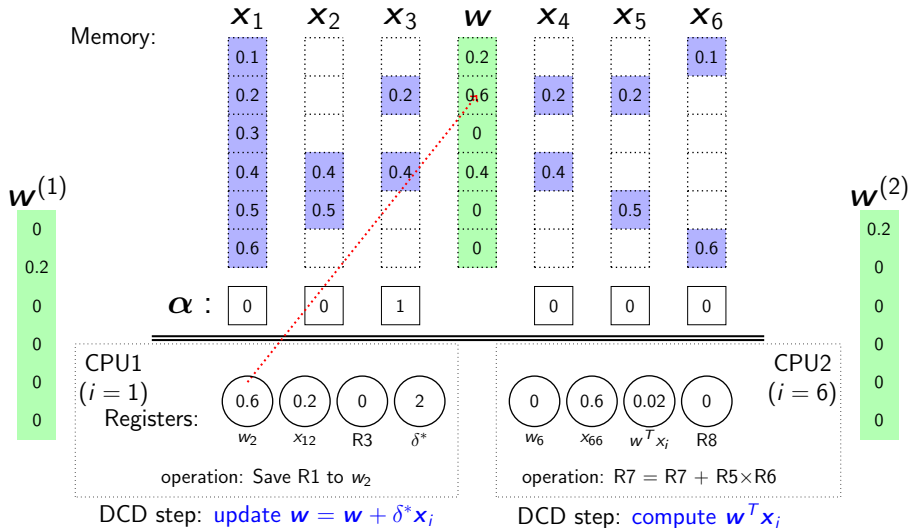
# Dual Coordinate Descent in Parallel

## Inconsistent Read of $w$

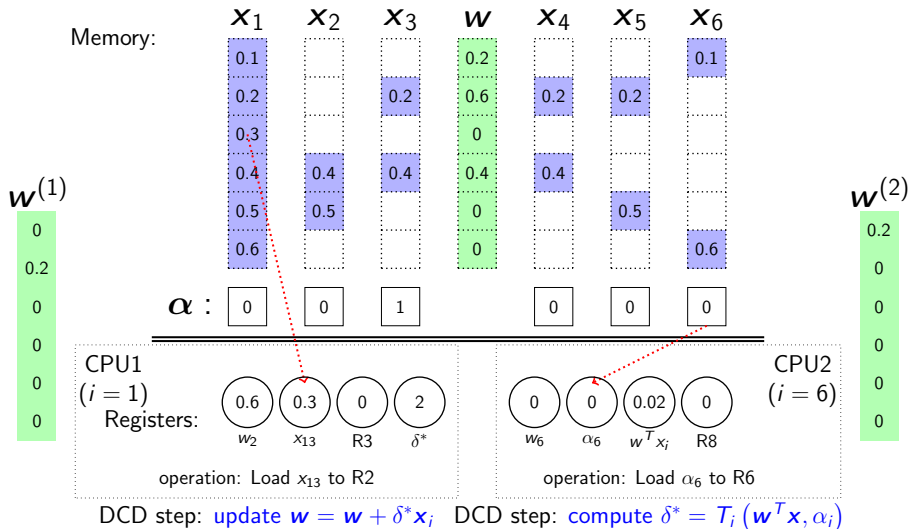


# Dual Coordinate Descent in Parallel

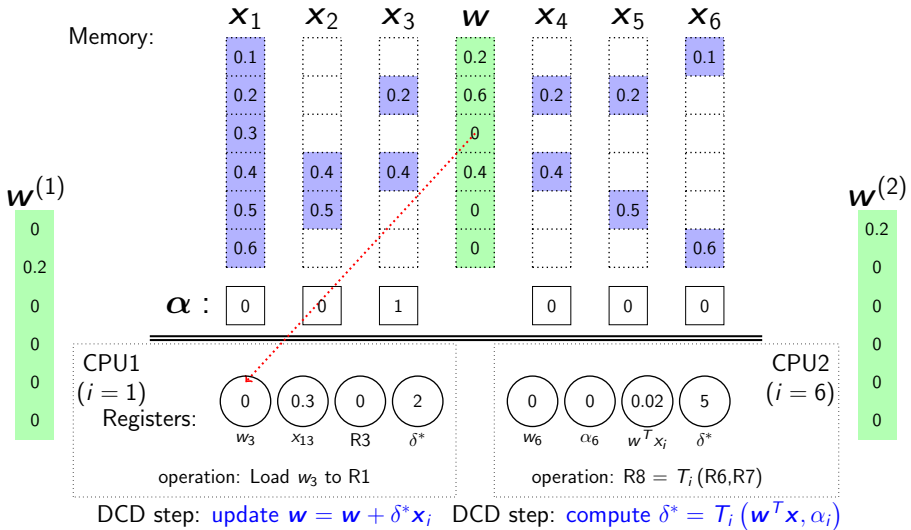
Inconsistent Read of  $w$



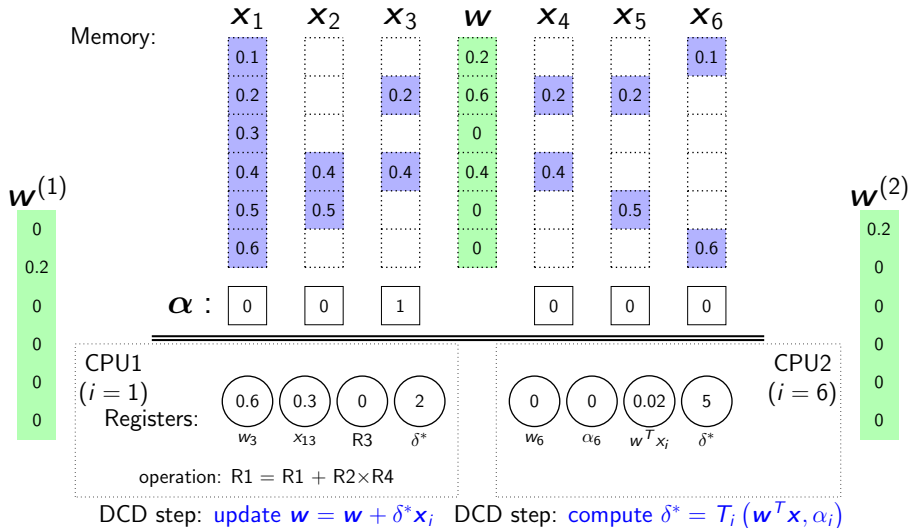
# Dual Coordinate Descent in Parallel



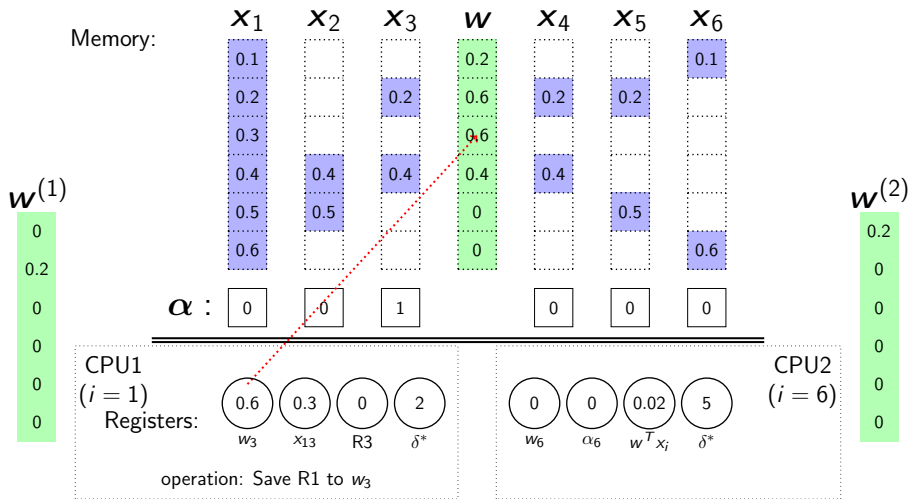
# Dual Coordinate Descent in Parallel



# Dual Coordinate Descent in Parallel



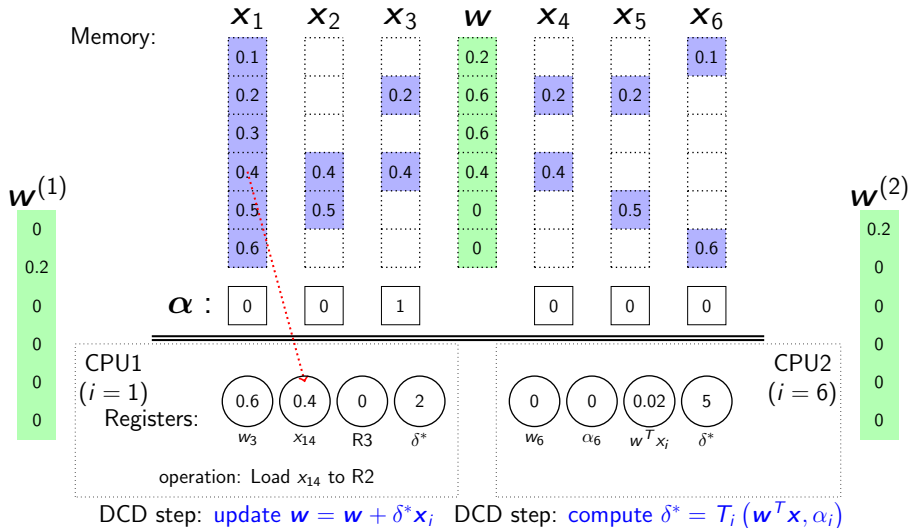
# Dual Coordinate Descent in Parallel



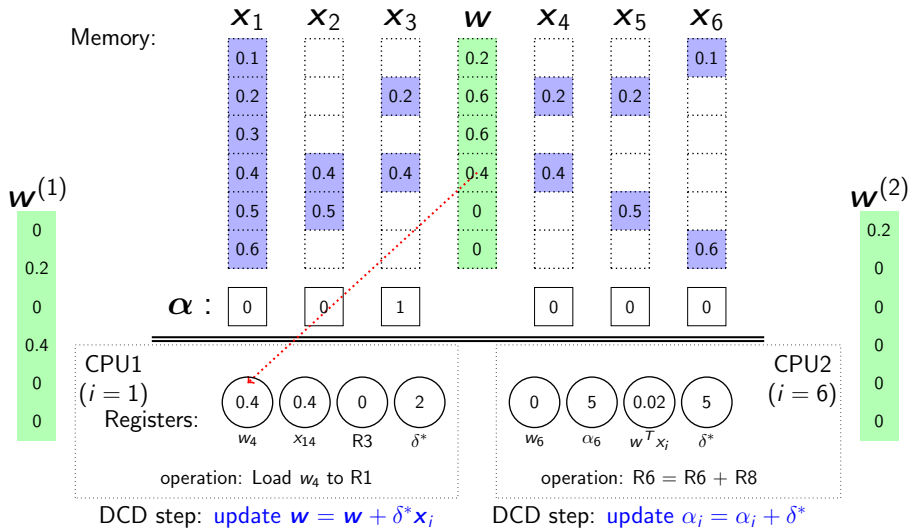
DCD step:  $\text{update } w = w + \delta^* x_i$    DCD step:  $\text{compute } \delta^* = T_i(w^T x, \alpha_i)$



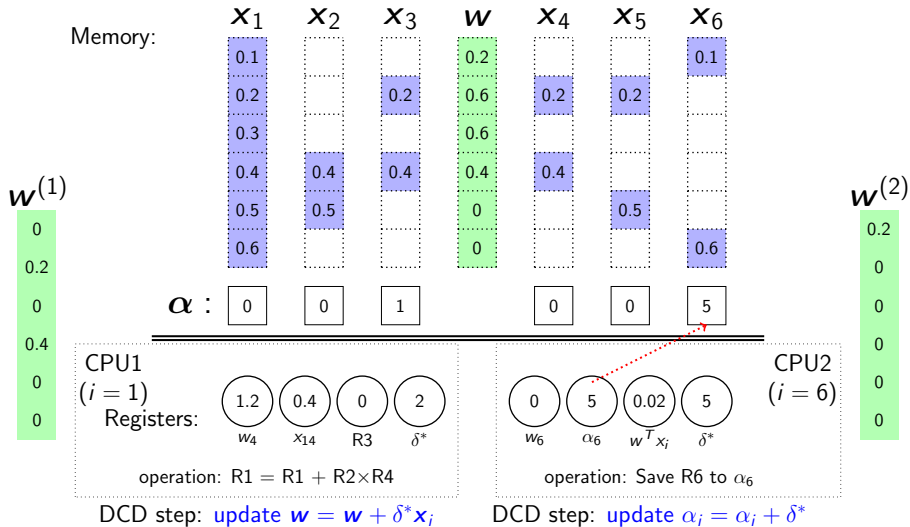
# Dual Coordinate Descent in Parallel



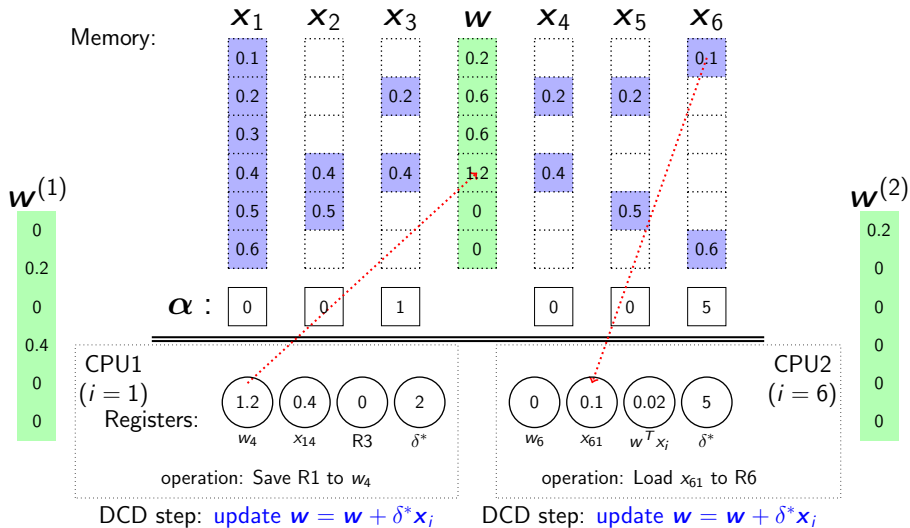
# Dual Coordinate Descent in Parallel



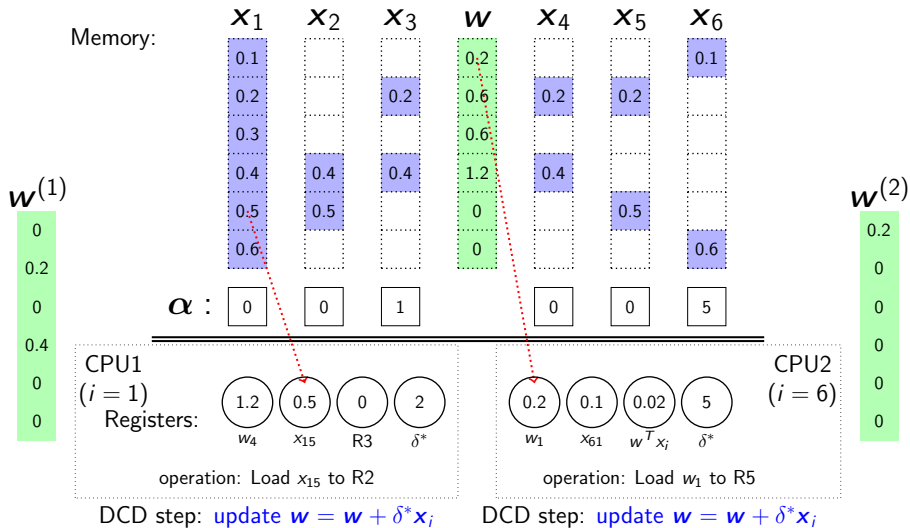
# Dual Coordinate Descent in Parallel



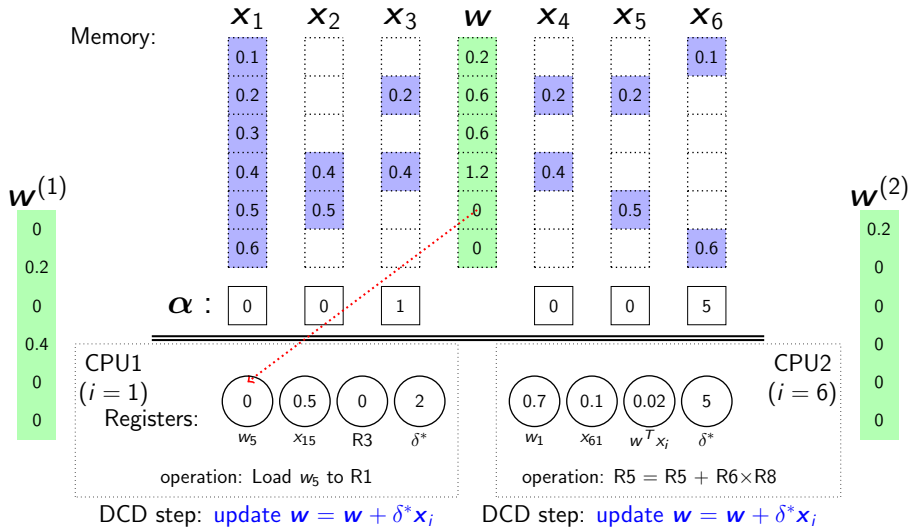
# Dual Coordinate Descent in Parallel



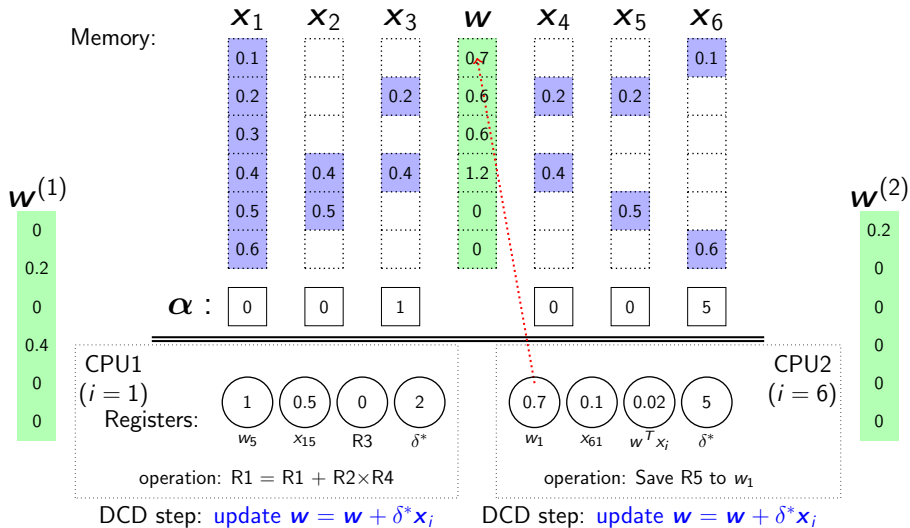
# Dual Coordinate Descent in Parallel



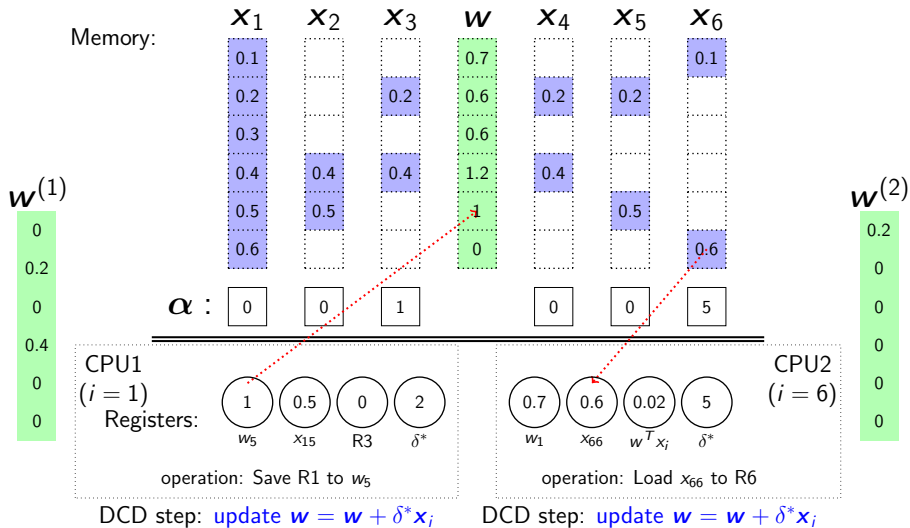
# Dual Coordinate Descent in Parallel



# Dual Coordinate Descent in Parallel

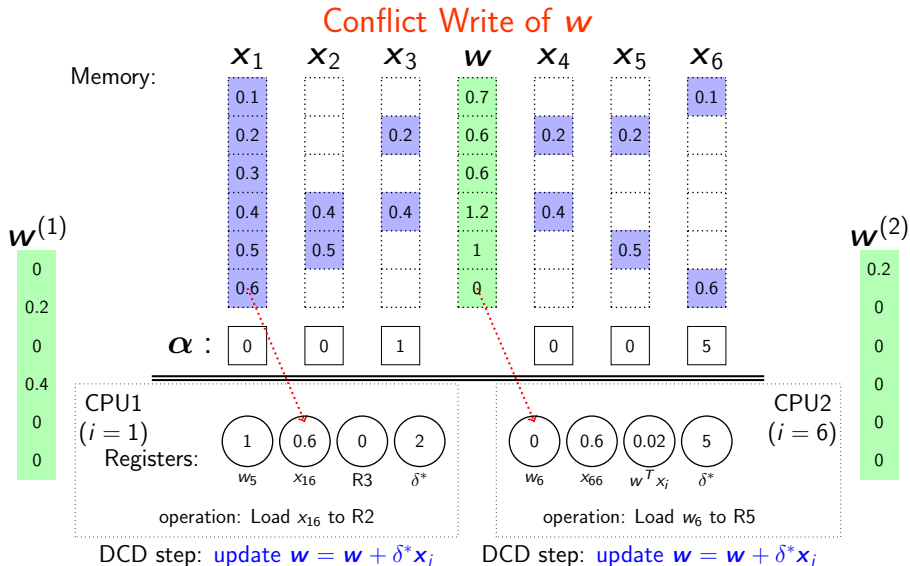


# Dual Coordinate Descent in Parallel

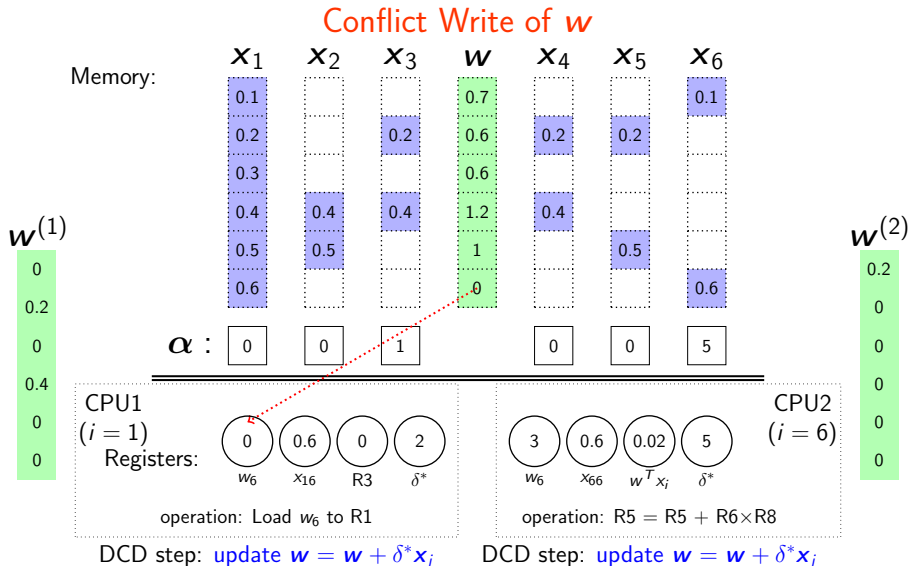




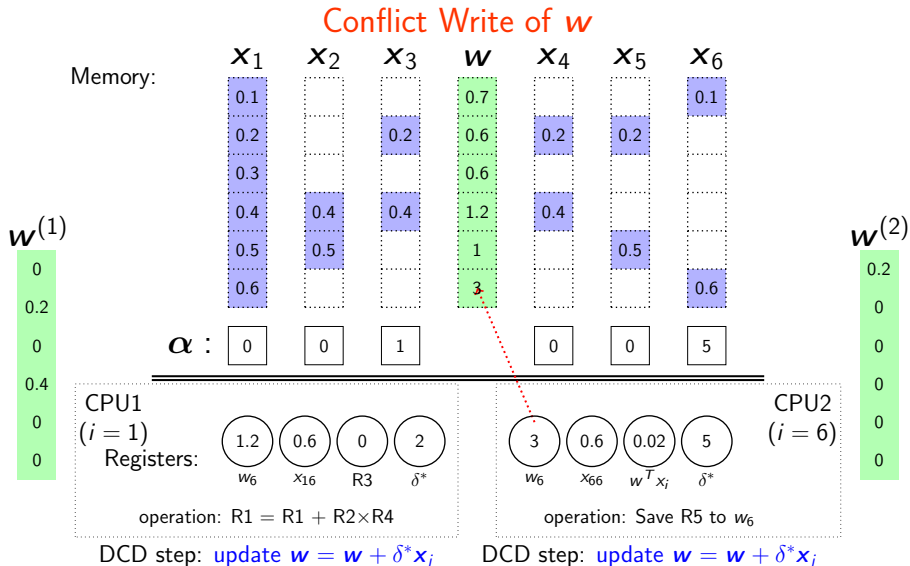
# Dual Coordinate Descent in Parallel



# Dual Coordinate Descent in Parallel

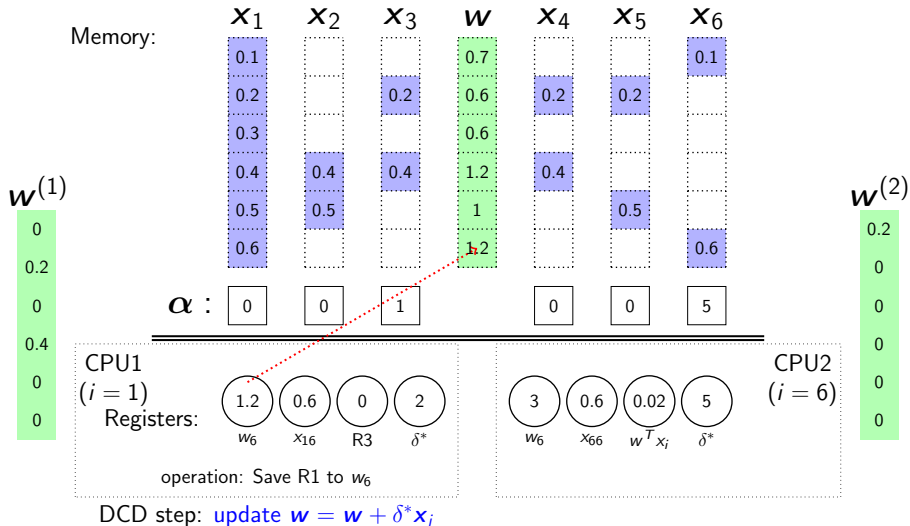


# Dual Coordinate Descent in Parallel

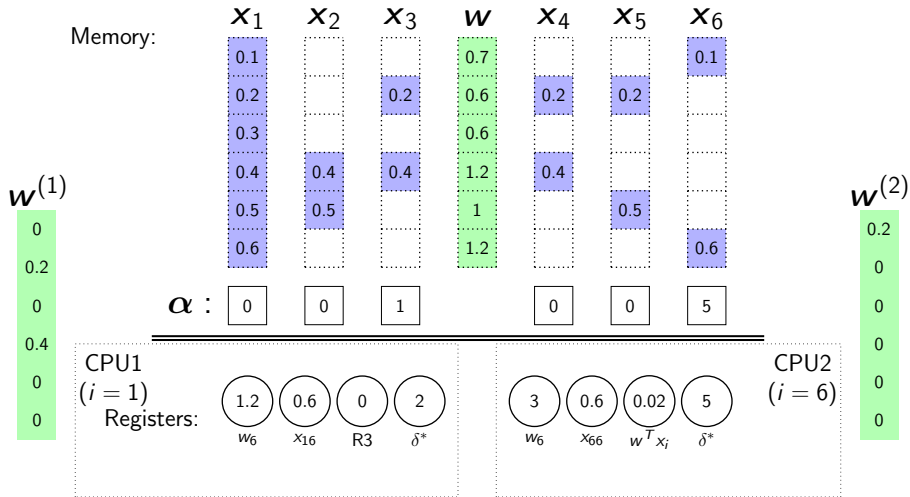


# Dual Coordinate Descent in Parallel

Conflict Write of  $w$



# Dual Coordinate Descent in Parallel



Each thread repeatedly performs the following updates.

For  $t = 1, 2, \dots$

- 1 Randomly pick an index  $i$
- 2 **Lock**  $\{w_j \mid (x_i)_j \neq 0\}$
- 3 Compute  $\mathbf{w}^T \mathbf{x}_i$
- 4 Update  $\alpha_i \leftarrow \alpha_i + \delta^*$  where  $\delta^* = T_i(\mathbf{w}^T \mathbf{x}_i, \alpha_i)$
- 5 Update  $\mathbf{w} \leftarrow \mathbf{w} + \delta^* \mathbf{x}_i$ .
- 6 **Unlock the variables.**

# How to Resolve the Issues

Three *PASSCoDe* approaches:

- **lock**: acquire **locks** for all necessary  $w_j$  before the update

	inconsistent read	conflict write
<b>PASSCoDe-Lock</b>	resolved	resolved

Scaling (on rcv1 with 100 epochs):

# threads	<b>Lock</b>
2	98.03s / 0.27x
4	106.11s / 0.25x
10	114.43s / 0.23x

Each thread repeatedly performs the following updates.

For  $t = 1, 2, \dots$

- 1 Randomly pick an index  $i$
- 2 Compute  $\mathbf{w}^T \mathbf{x}_i$
- 3 Update  $\alpha_i \leftarrow \alpha_i + \delta^*$  where  $\delta^* = T_i(\mathbf{w}^T \mathbf{x}_i, \alpha_i)$
- 4 For each  $j \in N(i)$ 
  - 5 Update  $w_j \leftarrow w_j + \delta^*(\mathbf{x}_i)_j$  **atomically**



# How to Resolve the Issues

Three *PASSCoDe* approaches:

- **lock**: acquire **locks** for all necessary  $w_j$  before the update
- **atomic**: apply **atomic** operation for  $w_j = w_j + \delta^* x_{ij}$

	inconsistent read	conflict write
<b>PASSCoDe-Lock</b>	resolved	resolved
<b>PASSCoDe-Atomic</b>	remained	resolved

Scaling (on rcv1 with 100 epochs):

# threads	Lock	Atomic
2	98.03s / 0.27x	15.28s / 1.75x
4	106.11s / 0.25x	8.35s / 3.20x
10	114.43s / 0.23x	3.86s / 6.91x

# Analysis for PASSCoDe-Atomic

- Atomic operations guarantee:
  - all updates to  $\mathbf{w}$  will be performed eventually
  - $\hat{\mathbf{w}} = \sum_{i=1}^n \hat{\alpha}_i \mathbf{x}_i$  holds for the outputted  $(\hat{\mathbf{w}}, \hat{\alpha})$
- **Bounded delay assumption**: to handle inconsistent read of  $\mathbf{w}$ 
  - all updates of  $\mathbf{w}$  before  $\tau$  iterations must be performed

## Theorem

Under certain conditions on  $\tau$ , **PASSCoDe – Atomic** has global linear convergence rate in expectation:

$$E [D(\alpha^{j+1}) - D(\alpha^*)] \leq \eta E [D(\alpha^j) - D(\alpha^*)]$$

**Our analysis covers logistic regression and SVM with hinge loss (where the dual problem is not strictly convex).**

Each thread repeatedly performs the following updates.

For  $t = 1, 2, \dots$

- 1 Randomly pick an index  $i$
- 2 Compute  $\mathbf{w}^T \mathbf{x}_i$
- 3 Update  $\alpha_i \leftarrow \alpha_i + \delta^*$  where  $\delta^* = T_i(\mathbf{w}^T \mathbf{x}_i, \alpha_i)$
- 4 Update  $\mathbf{w} \leftarrow \mathbf{w} + \delta^* \mathbf{x}_i$

# How to Resolve the Issues

Three *PASSCoDe* approaches:

- **lock**: acquire **locks** for all necessary  $w_j$  before the update
- **atomic**: apply **atomic** operation for  $w_j = w_j + \delta^* x_{ij}$
- **wild**: do **nothing** to resolve either issue

	inconsistent read	conflict write
<b>PASSCoDe-Lock</b>	resolved	resolved
<b>PASSCoDe-Atomic</b>	remained	resolved
<b>PASSCoDe-Wild</b>	remained	remained

Scaling (on rcv1 with 100 epochs):

# threads	Lock	Atomic	Wild
2	98.03s / 0.27x	15.28s / 1.75x	14.08s / 1.90x
4	106.11s / 0.25x	8.35s / 3.20x	7.61s / 3.50x
10	114.43s / 0.23x	3.86s / 6.91x	3.59s / 7.43x

# Analysis for PASSCoDe-Wild

- Some updates are missing due to memory conflicts

- Which one for prediction,  $\hat{\mathbf{w}}$  or  $\bar{\mathbf{w}}$ ?

- for the final ( $\hat{\mathbf{w}}$ ,  $\hat{\alpha}$ ):

$$\hat{\mathbf{w}} \neq \sum_{i=1}^n \hat{\alpha}_i \mathbf{x}_i$$

- construct  $\bar{\mathbf{w}}$  from the final  $\hat{\alpha}$ :

$$\bar{\mathbf{w}} = \sum_{i=1}^n \hat{\alpha}_i \mathbf{x}_i$$

	# threads	Prediction Accuracy (%) by		
		$\hat{\mathbf{w}}$	$\bar{\mathbf{w}}$	LIBLINEAR
news20	4	97.1	96.1	97.1
	8	97.2	93.3	
covtype	4	67.8	38.0	66.3
	8	67.6	38.0	
rcv1	4	97.7	97.5	97.7
	8	97.7	97.4	
webspam	4	99.1	93.1	99.1
	8	99.1	88.4	
kddb	4	88.8	79.7	88.8
	8	88.8	87.7	

Question: why  $\hat{\mathbf{w}}$  is better than  $\bar{\mathbf{w}}$ ?

# Backward Analysis for PASSCoDe-Wild

Recall the primal problem

$$\mathbf{w}^* = \arg \min_{\mathbf{w}} P(\mathbf{w}) := \frac{1}{2} \|\mathbf{w}\|^2 + \sum_{i=1}^n \ell_i(\mathbf{w}^T \mathbf{x}_i)$$

## Theorem

Let  $\epsilon$  be the error caused by the memory conflicts.

$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w}} \hat{P}(\mathbf{w}) := \frac{1}{2} \|\mathbf{w} + \epsilon\|^2 + \sum_{i=1}^n \ell_i(\mathbf{w}^T \mathbf{x}_i)$$

$$\bar{\mathbf{w}} = \arg \min_{\mathbf{w}} \bar{P}(\mathbf{w}) := \frac{1}{2} \|\mathbf{w}\|^2 + \sum_{i=1}^n \ell_i((\mathbf{w} - \epsilon)^T \mathbf{x}_i)$$

- $\hat{P}(\mathbf{w})$  is the problem with the perturbation on the regularization term
- $\bar{P}(\mathbf{w})$  is the problem with the perturbation on the prediction term

# Datasets and Experimental Settings

## Datasets.

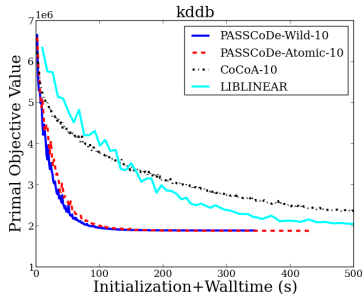
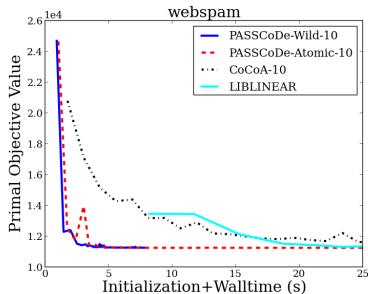
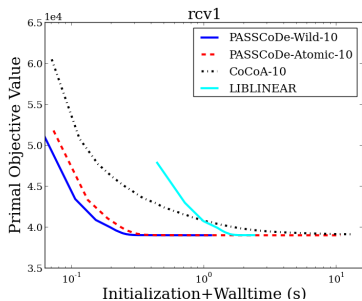
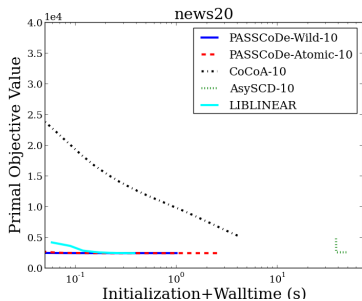
	$n$	$\tilde{n}$	$d$	$\bar{d}$	$C$
news20	16,000	3,996	1,355,191	455.5	2
rcv1	677,399	20,242	47,236	73.2	1
webspam	280,000	70,000	16,609,143	3727.7	1
kddb	19,264,097	748,401	29,890,095	29.4	1

## Compared Implementation.

- LIBLINEAR: serial baseline
- *PASSCoDe-Wild* and *PASSCoDe-Atomic*: our methods
- *CoCoA*: a multi-core version of [Jaggi et al, 2014]
- *AsySCD*: [Liu & Wright, 2014]

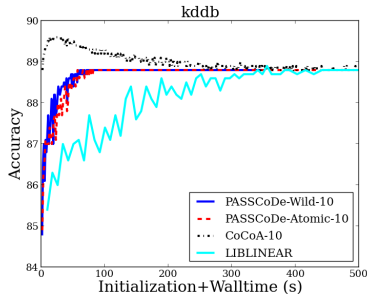
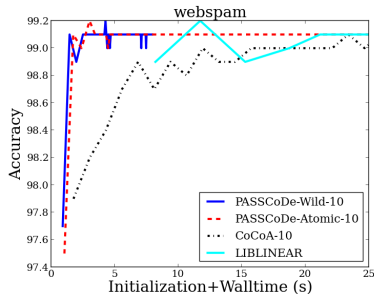
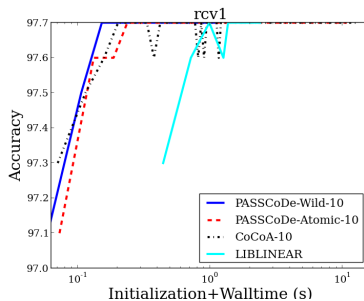
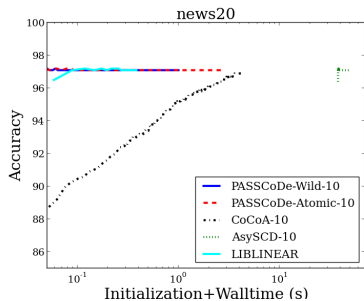
**Machine:** Intel Multi-core machine with 256 GB Memory

# Convergence in terms of Waltime

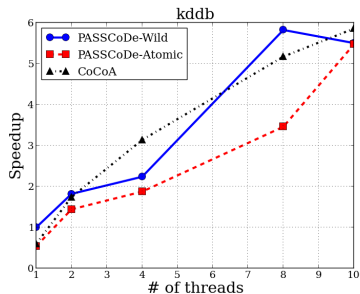
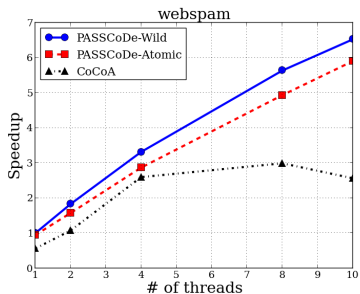
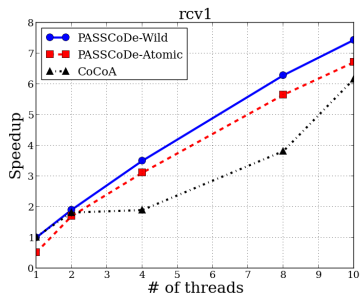
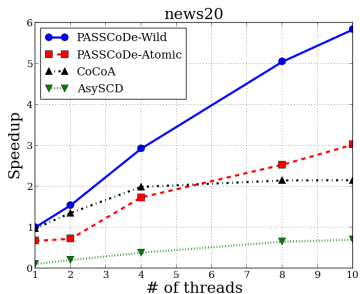




# Accuracy



# Speedup



# Conclusions

- *PASSCoDe*: an simple but effective asynchronous dual coordinate descent
- Analysis three variants
  - **PASSCoDe-Lock**
  - **PASSCoDe-Atomic**: established global linear convergence
  - **PASSCoDe-Wild**: backward analysis
- Future work: extend the analysis to L1-regularized problems
  - LASSO
  - L1-regularized Logistic Regression