

Timestamp of the last model used by worker p

- $l_{\mathcal{D}}$ Model (parameter vector) at epoch t $\boldsymbol{x}_t$ Learning rate at epoch t $\gamma_t$
- Each gradient is a tuple  $[\boldsymbol{g}_p, l]$  denoting that  $g_p$ a worker p computed the gradient  $\boldsymbol{g}_p$  w.r.t  $\boldsymbol{x}_l$
- Set of gradients that the server receives in epoch t
- Staleness value for a gradient  $[\boldsymbol{g}, l]$  at epoch t  $\tau_{tl}$  $(\tau_{tl} \triangleq k - l)$
- Mini-batch of training examples

## 4. Kardam

• Byzantine resilience against f/n workers, f <= n/3

$$\frac{n-2f}{n-f} \le SL \le \frac{n-f}{n}$$

$$\boldsymbol{x}_{t+1} = \boldsymbol{x}_t - \gamma_t \sum_{[\boldsymbol{G}(\boldsymbol{x}_l; \boldsymbol{\xi}_m), l] \in \boldsymbol{\mathcal{G}}_t} \Lambda(\tau_{tl}) \cdot \boldsymbol{G}(\boldsymbol{x}_l; \boldsymbol{\xi}_m)$$

#### => convergence rate bound

# 7. Evaluation

## <u>Setup</u>

- CIEAD 100	Parameters	Input	Conv1	Pool1	Conv2	Pool2	FC1	FC2	FC3
<ul> <li>CIFAR-100</li> <li>CNN</li> </ul>	Kernel size	$32 \times 32 \times 3$	$3 \times 3 \times 16$	$3 \times 3$	$3 \times 3 \times 64$	$4 \times 4$	384	192	100
	Strides			3×3		4×4			

• f = 3, n = 10 workers

 Baseline-ASGD: no dampening component • SSGD: ideal (synchronous) SGD •  $\Lambda_1 = 1/(\tau+1), \Lambda_2 = \exp(0.5\tau), \Lambda_3 = \exp(0.2\tau)$ 

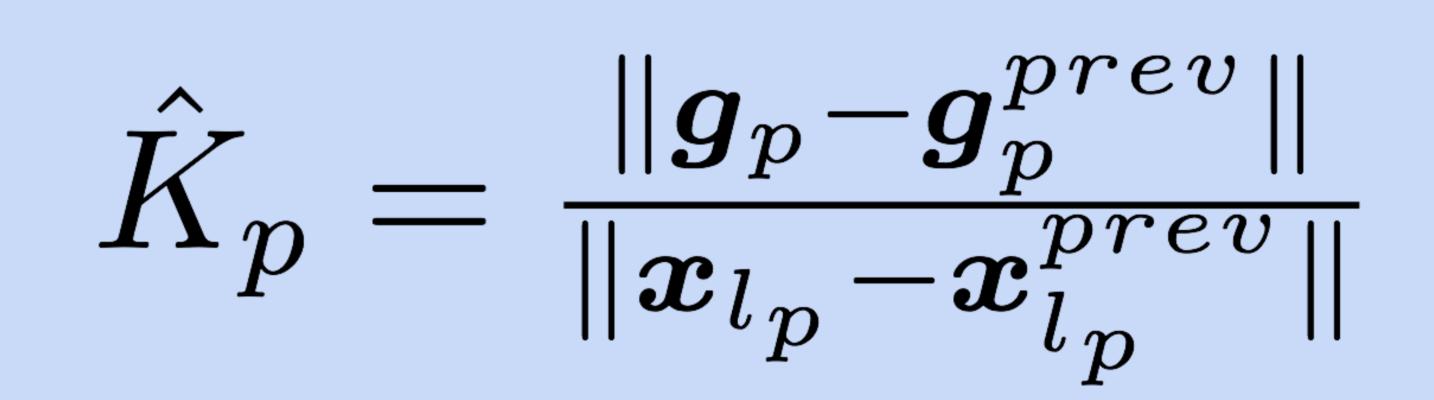
#### Provable (almost sure) convergence

# 5. Filtering component

a. Lipschitz filter empirical Lipschitz coefficient at worker p:

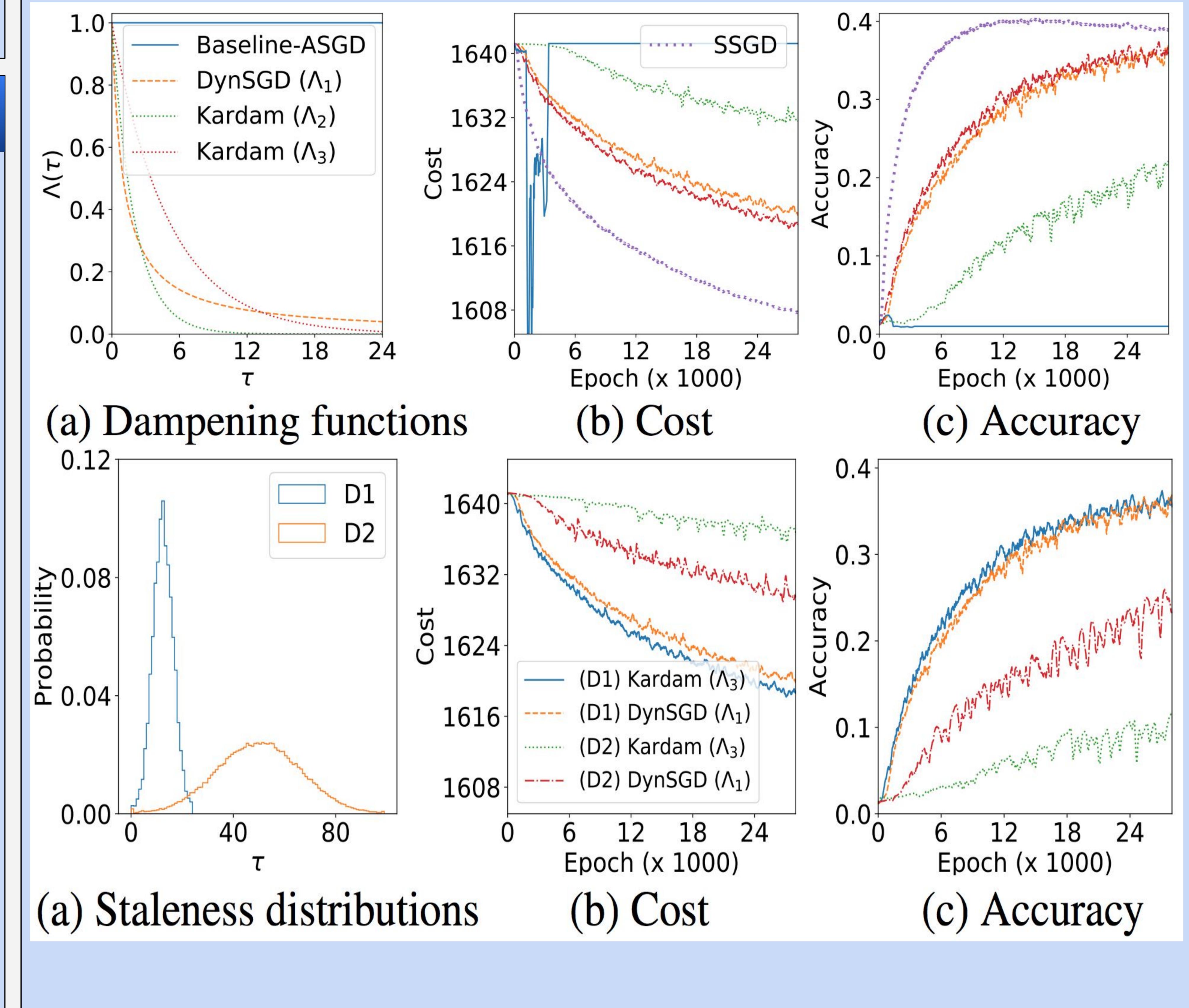
empirical Lipschitz coefficient at server after response from p:

 $\sim$  $\boldsymbol{\wedge}$  $\tilde{K}_t^p \leq K_t \triangleq quantile_{\underline{n-f}} \{K_p\}_{p \in P}$ 



 $\|\boldsymbol{g}_p - \boldsymbol{g}_q\|$ 

 $\| \boldsymbol{x}_t - \boldsymbol{x}_{t-1} \|$ 



### b. Frequency filter

limits the number of successive gradients from a single worker to a value of f

#### Slowdown

a ^ b => correct cone:  $\langle \mathbb{E}_{\xi} G(x;\xi), \nabla Q(x) \rangle > \Omega((\|\nabla Q(x_t)\| - \sqrt{d}\sigma) \|\nabla Q(x_t)\|)$   $\langle \mathbb{E}_{\xi} G(x;\xi), \nabla Q(x) \rangle > \Omega((\|\nabla Q(x_t)\| - \sqrt{d}\sigma) \|\nabla Q(x_t)\|)$   $\Lambda_1: 27.9 \% \text{ filtered gradients } (D_1)$   $\Lambda_3: 19.6 \% \text{ filtered gradients } (D_1)$ 

# **Sake away**

 gradient filtering => Byzantine resilience aradient dampening => Asynchronous convergence

International Conference on Machine Learning (ICML) 2018, Stockholm, Sweden