

# A Delay-tolerant Proximal-Gradient Algorithm for Distributed Learning

K. Mishchenko<sup>1</sup>, F. Lutzeler<sup>2</sup>, J. Malick<sup>2</sup>, and M.-R. Amini<sup>2</sup>



<sup>1</sup>KAUST <sup>2</sup>Univ. Grenoble Alpes <sup>3</sup>CNRS



## Abstract

### Setting:

- **Distributed** Data distributed across  $M$  slaves
- **Asynchronous** Slaves send their results irregularly
- **Regularized** Regularized by a non-smooth function, such as  $\ell_1$ -penalty
- **Unbalanced** Computational power of slaves can be very uneven
- **Scarce communications** Transmission cost can be very high

### Results:

- We propose an **algorithm** based on **gradient** and **proximal** steps
- For strongly convex and smooth functions, rate is **linear** even when **delays are unbounded**

### Applications:

- $\ell_1$  regularized logistic regression, group lasso minibatch processing, etc.

## PROBLEM FORMULATION

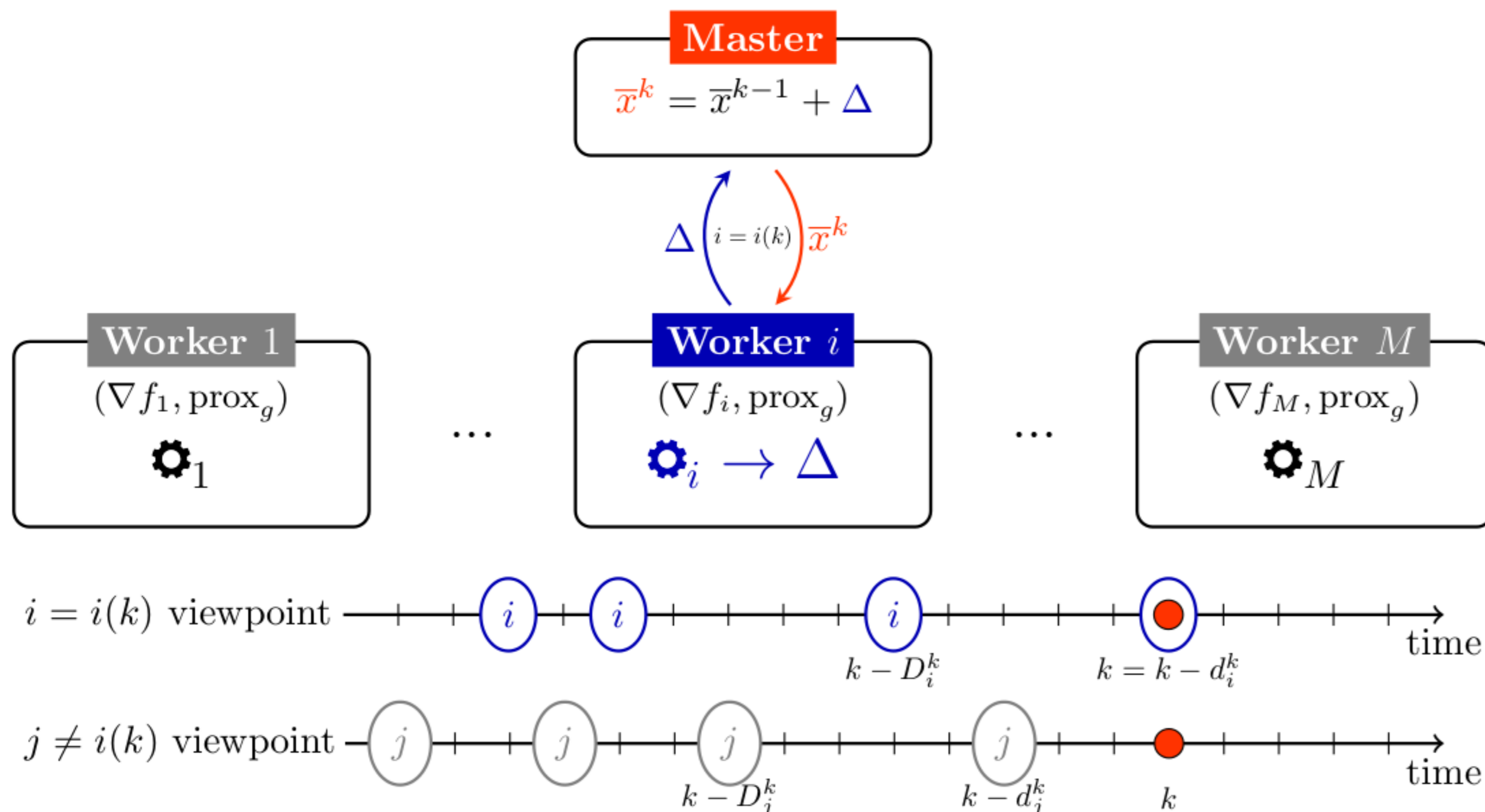
$$\min_{x \in \mathbb{R}^d} \frac{1}{n} \sum_{i=1}^n \ell_i(x) + r(x),$$

Global Objective  $\downarrow$  Local Data

$$\min_{x \in \mathbb{R}^d} \sum_{i=1}^M \underbrace{\left( \frac{n_i}{n} \right)}_{\text{proportion}} \underbrace{\left( \frac{1}{n_i} \sum_{j \in \mathcal{S}_i} \ell_j(x) \right)}_{\text{local empirical loss}} + r(x)$$

- $\mathcal{S}_i$  local dataset at  $i$
- Each  $f_i$  is convex,  $L$ -Smooth  
 $\|\nabla f_i(x) - \nabla f_i(y)\| \leq L \|x - y\|$
- $\mu$ -strongly convex
- $r$  is convex

## ASYNCHRONOUS FRAMEWORK



- **iteration** = receive from a worker + master update + send back
- **time  $k$**  = number of iterations
- **delay  $d_i^k$**  = time since last exchange with  $i$   
 $d_i^k = 0$  iff  $i$  updates at time  $k$ ,  $d_i^k = d_i^{k-1} + 1$  elsewhere
- **second delay  $D_i^k$**  = time since penultimate exchange with  $i$

## DAVE COMMUNICATION SCHEME

$$\bar{x}^k = \bar{x}^{k-1} + \Delta \text{ with } \Delta = \pi_i (x_i^k - x_i^{k-D_i^k}) \text{ for } i = i(k)$$

$$\text{i.e. } \bar{x}^k = \sum_{i=1}^M \pi_i x_i^{k-d_i^k} = \sum_{i=1}^M \pi_i \text{gear}_i (x_i^{k-D_i^k})$$

- master variable  $\bar{x}^k$  = combination of workers last contributions  $(x_i^{k-d_i^k})_i$   
 $\pi_i$  = proportion of data at worker  $i$
- one update/time = one worker contribution **but** all workers are always involved in the master variable with even proportions

## REVISITING THE CLOCK

- **epoch sequence  $(k_m)$**  recursively defined by  $k_0 = 0$  and  
 $k_{m+1} = \min\{k : \text{each worker made at least 2 updates on the interval } [k_m, k]\}$   
 $= \min\{k : k - D_i^k \geq k_m \text{ for all } i = 1, \dots, M\}$
- **epoch time  $m$**  = number of epochs

## DAVE-RPG

Worker  $i$  performs **Repeated Proximal Gradient** steps

### Master:

```
Initialize  $\bar{x}$ 
while not converged do
  when a worker finishes:
    Receive adjustment  $\Delta$  from it
     $\bar{x} \leftarrow \bar{x} + \Delta$ 
    Send  $\bar{x}$  to the agent in return
   $k \leftarrow k + 1$ 
end
Interrupt all slaves
Output  $x = \text{prox}_{\gamma g}(\bar{x})$ 
```

### Worker $i$ :

```
Initialize  $x = x_i = \bar{x}$ ,
while not interrupted by master do
  Receive the most recent  $\bar{x}$ 
  Select a number of repetitions  $p$ 
  Initialize  $\Delta = 0$ 
  for  $q = 1$  to  $p$  do
     $z \leftarrow \text{prox}_{\gamma g}(\bar{x} + \Delta)$ 
     $x_i \leftarrow z - \gamma \nabla f_i(z)$ 
     $\Delta \leftarrow \Delta + \pi_i (x_i - x_i^{\text{prev}})$ 
     $x_i \leftarrow x_i^{\text{prev}}$ 
  end
  Send the adjustment  $\Delta$  to the master
end
```

## CONVERGENCE RESULTS

Take  $p = 1$ . Then, for  $\gamma \in (0, 2/(\mu + L)]$ ,

$$\forall k \geq k_m, \quad \|x^k - x^*\|^2 \leq \left(1 - \frac{2\gamma\mu L}{\mu + L}\right)^m \|x^0 - x^*\|^2$$

where  $x^*$  is the unique minimizer of  $\min_x \sum_{i=1}^M \pi_i f_i(x) + g(x)$ .

- Usual fixed stepsize, delay-independent
- Linear rate along the *epoch* sequence; delay-tolerant

Take  $p$  changing across time and workers. Then, for  $\gamma \in (0, 2/(\mu + L)]$ ,

$$\forall k \geq k_m, \quad \|x^k - x^*\|^2 \leq \left(1 - \frac{2\gamma\mu L}{\mu + L}\right)^m \prod_{\ell=1}^m \alpha_\ell \|x^0 - x^*\|^2,$$

with epoch improvement

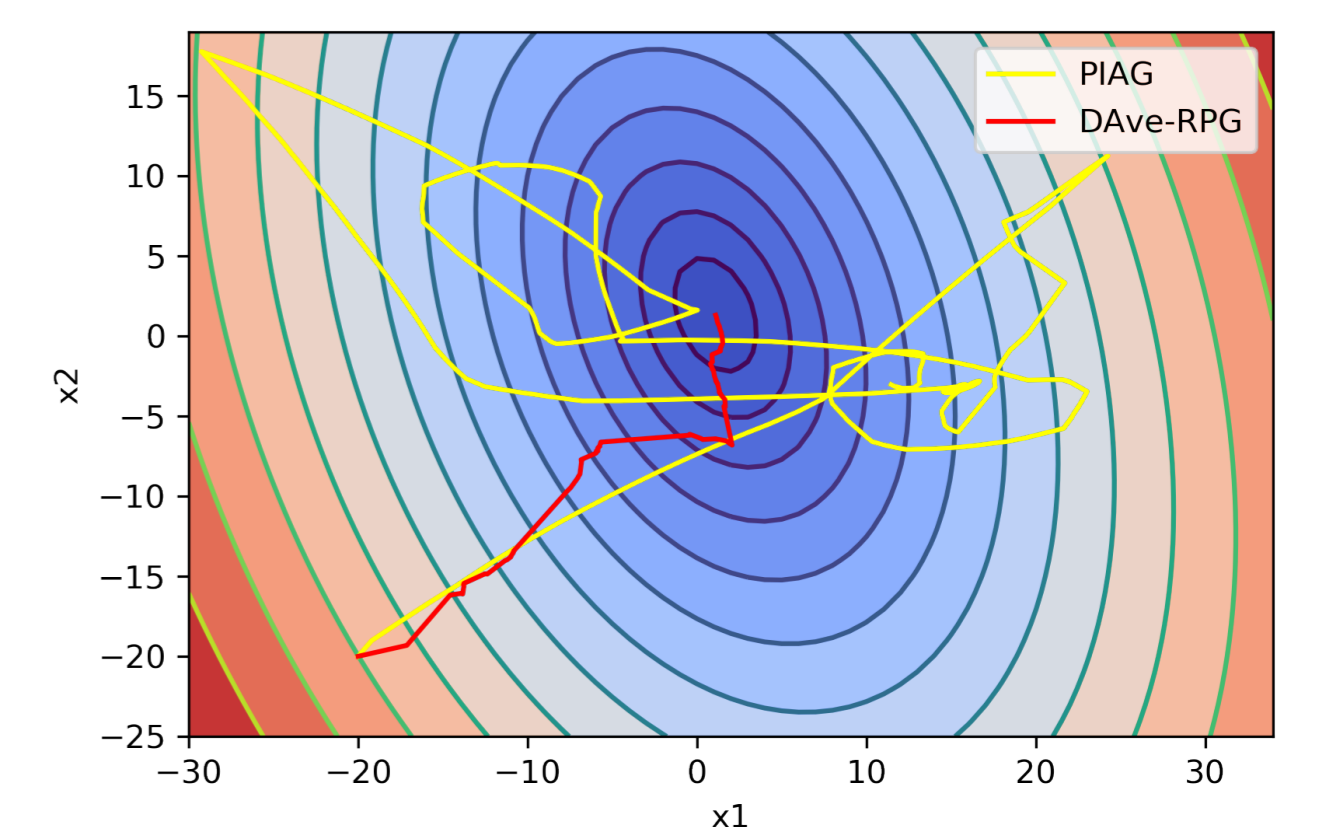
$$\alpha_\ell = \max_{i, k \in [k_\ell, k_{\ell+1})} 1 - \gamma\mu \sum_{q=1}^{p_i^k - 1} (1 - \gamma\mu)^{q-1} \min_i \pi_i^q.$$

- Workers compute more without coordinating
- Compromise Communication/Computation

## COMPARISON WITH GRADIENT COMBINATION

	DAVE-RPG ( $p = 1$ ) iterates	PIAG gradients
Combining:		
Update $x^k$ :	$\text{prox}_{\gamma g} \left( \sum_{i=1}^M \pi_i x_i^{k-D_i^k} - \gamma \sum_{i=1}^M \pi_i \nabla f_i(x_i^{k-D_i^k}) \right)$	$\text{prox}_{\gamma g} \left( x^{k-1} - \gamma \sum_{i=1}^M \pi_i \nabla f_i(x_i^{k-D_i^k}) \right)$

- Combining *iterates* is more stable than combining *gradients*
- Example: 2D quadratic functions on 5 worker but one worker 10x slower  
Stepsize  $\gamma$  of PIAG is 10x smaller due to delays while the one for DAVE-PG stays =



## NUMERICAL RESULTS

- Logistic regression w/ elastic net  $\frac{1}{m} \sum_{j=1}^m \log(1 + \exp(-y_j z_j^T x)) + \lambda_1 \|x\|_1 + \frac{\lambda_2}{2} \|x\|_2^2$
- 100 machines (1 CPU, 1 GB) in a cluster 10% of the data on machine 1, even on the rest

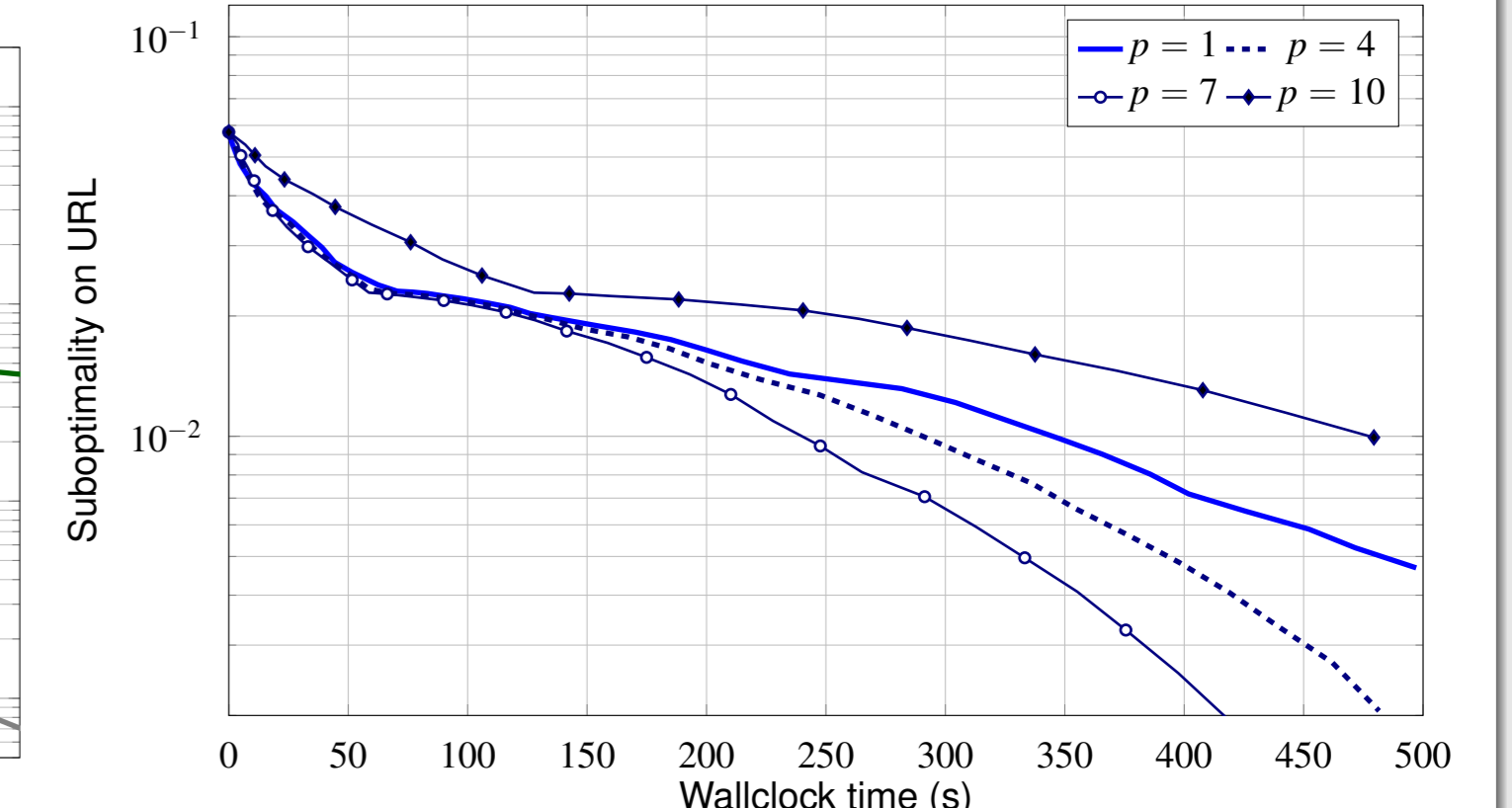
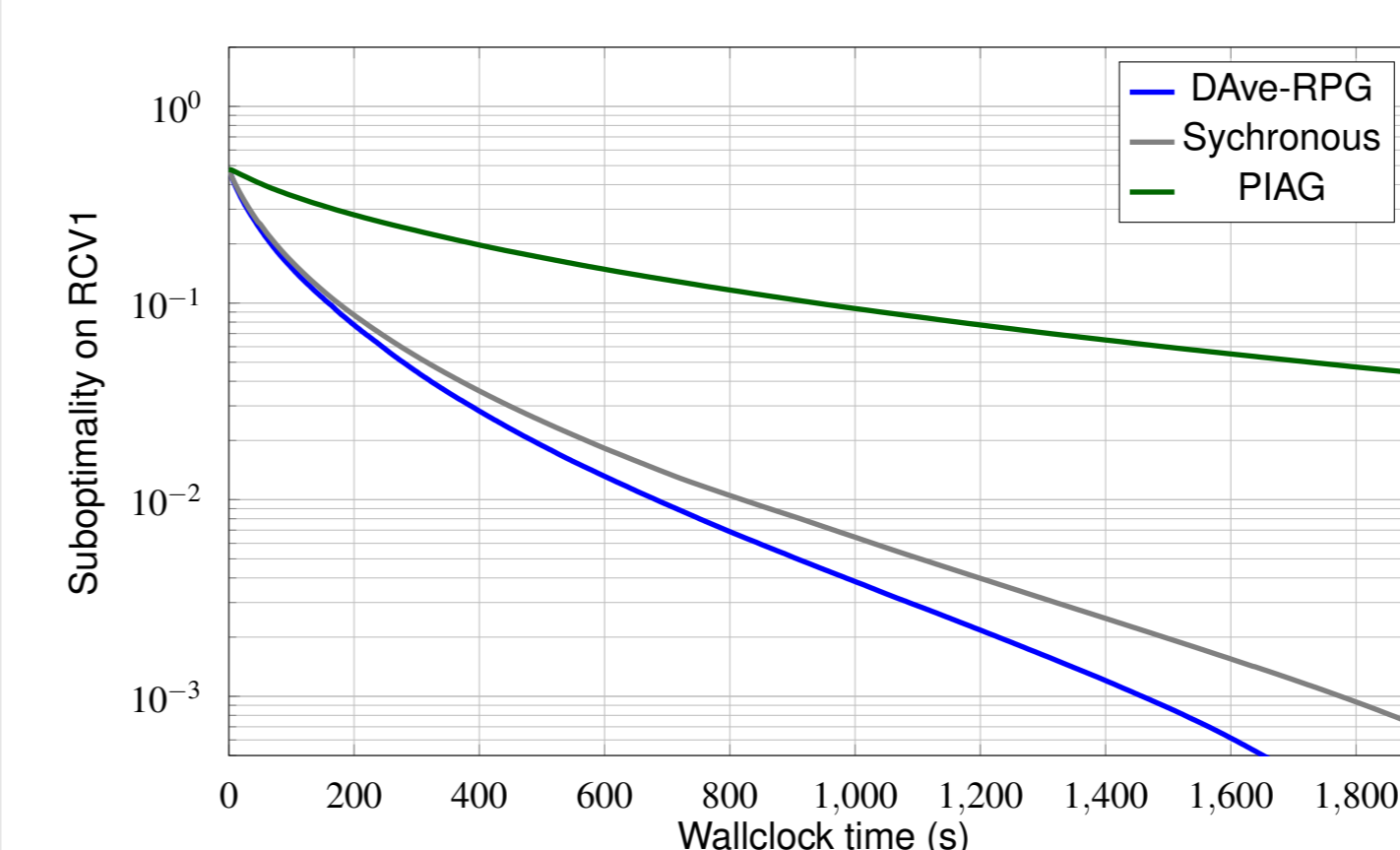


Illustration of delays on 10 evenly loaded machines

