# GIANT: Experiments

# Experiment Environment

- Spark 2.1.1     +     Scala 2.11.8

# Experiment Environment

- Spark 2.1.1    +    Scala 2.11.8
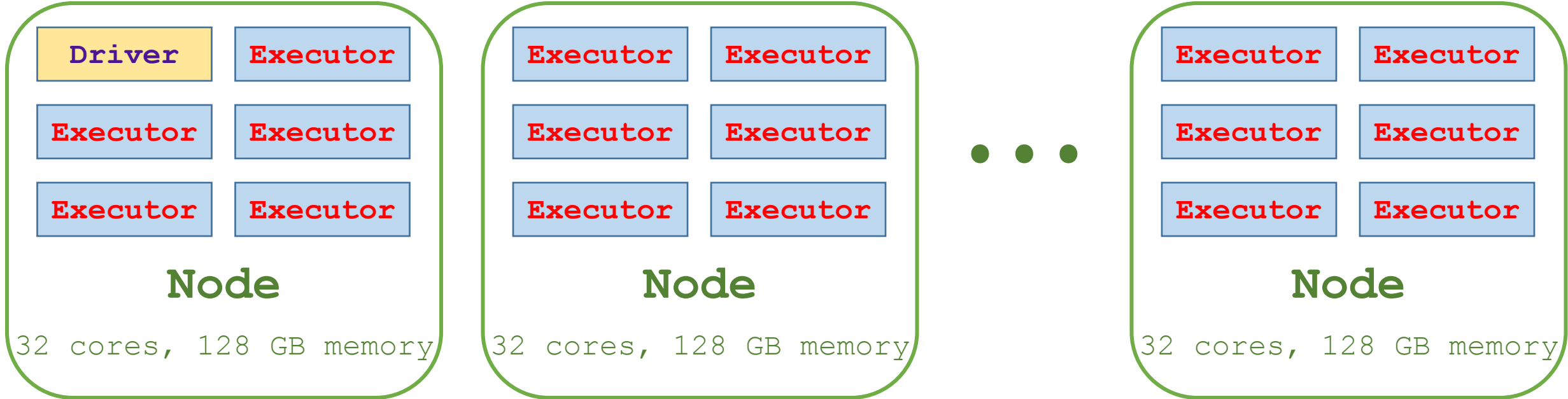


- Cori Supercomputer

# Experiment Environment

| Node | Node | | Node |
|---|---|---|---|
| **Driver** **Executor** | **Executor** **Executor** | | **Executor** **Executor** |
| **Executor** **Executor** | **Executor** **Executor** | $\cdots$ | **Executor** **Executor** |
| **Executor** **Executor** | **Executor** **Executor** | | **Executor** **Executor** |
| 32 cores, 128 GB memory | 32 cores, 128 GB memory | | 32 cores, 128 GB memory |

10 nodes, 59 executors

# Settings

- Solve the $\ell_2$-regularized logistic regression:

$$\min_{\mathbf{w}\in\mathbb{R}^d} \left\{ f(\mathbf{w}) \triangleq \frac{1}{n}\sum_{j=1}^{n} \log\left(1 + e^{-y_j \mathbf{x}_j^T \mathbf{w}}\right) + \frac{\gamma}{2}\|\mathbf{w}\|_2^2 \right\}$$

- Split $\mathbf{X} \in \mathbb{R}^{n \times d}$ (by data) to $m = 59$ parts.

- Local sample size $s = \dfrac{n}{m}$

# Settings

- Solve the $\ell_2$-regularized logistic regression:

$$\min_{\mathbf{w} \in \mathbb{R}^d} \left\{ f(\mathbf{w}) \triangleq \frac{1}{n} \sum_{j=1}^{n} \log\left(1 + e^{-y_j \mathbf{x}_j^T \mathbf{w}}\right) + \frac{\gamma}{2} \|\mathbf{w}\|_2^2 \right\}$$

- Split $\mathbf{X} \in \mathbb{R}^{n \times d}$ (by data) to $m = 59$ parts.

- Local sample size $s = \dfrac{n}{m} \quad \gtrsim \quad$ number of features $d$.

# Settings

- Solve the $\ell_2$-regularized logistic regression:

$$\min_{\mathbf{w} \in \mathbb{R}^d} \left\{ f(\mathbf{w}) \triangleq \frac{1}{n} \sum_{j=1}^{n} \log\left(1 + e^{-y_j \mathbf{x}_j^T \mathbf{w}}\right) + \frac{\gamma}{2} \|\mathbf{w}\|_2^2 \right\}$$

- Split $\mathbf{X} \in \mathbb{R}^{n \times d}$ (by data) to $m = 59$ parts.

- Local sample size $s = \dfrac{n}{m} \quad \gtrsim \quad$ number of features $d$.

- We use dense $\mathbf{X}$.

# Compared Methods

- Accelerated gradient descent (AGD)
    - choose *step size* from {0.1, 1, 10, 100}
    - choose *momentum* from {0.5, 0.9, 0.95, 0.99, 0.999}

# Compared Methods

- Accelerated gradient descent (AGD)

- Limited memory BFGS (a quasi-Newton method)
  - choose *number of history* from {30, 100, 300}
  - line search is used

# Compared Methods

- Accelerated gradient descent (AGD)

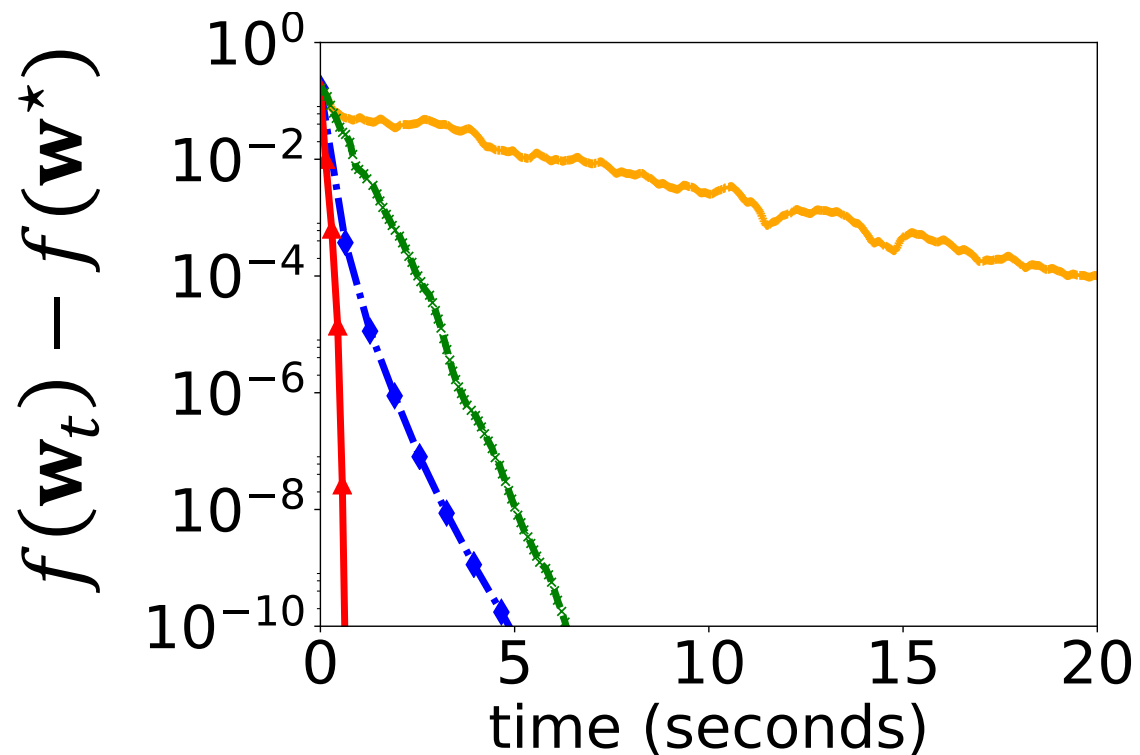- Limited memory BFGS

- DANE (another Newton-type method) [Shamir et al. 2014]

    - local solver: SVRG

    - choose *step size of SVRG* from {0.1, 1, 10, 100}

    - choose *max iteration of SVRG* from {30, 100, 300}

**Reference:**
Shamir, Srebro, & Zhang. Communication Efficient Distributed Optimization using an Approximate Newton-type Method. In *ICML*, 2014.

# Compared Methods

- Accelerated gradient descent (AGD)

- Limited memory BFGS

- DANE (another Newton-type method)

- GIANT
  - local solver: conjugate gradient (CG)
  - choose *max iteration of CG* from {30, 100, 300}

# Compared Methods

- Accelerated gradient descent (AGD)

- Limited memory BFGS

- DANE (another Newton-type method)

- GIANT

2 Tuning Parameters

1 Tuning Parameter

2 Tuning Parameters

1 Tuning Parameter

# Covtype (n=581K, d=54)

$\gamma = 10^{-4}$

$\gamma = 10^{-6}$



Legend: ADMM, AGD, DANE, GIANT, L-BFGS, AGD

# Epsilon (n=400K, d=2K)

$\gamma = 10^{-4}$

$\gamma = 10^{-6}$



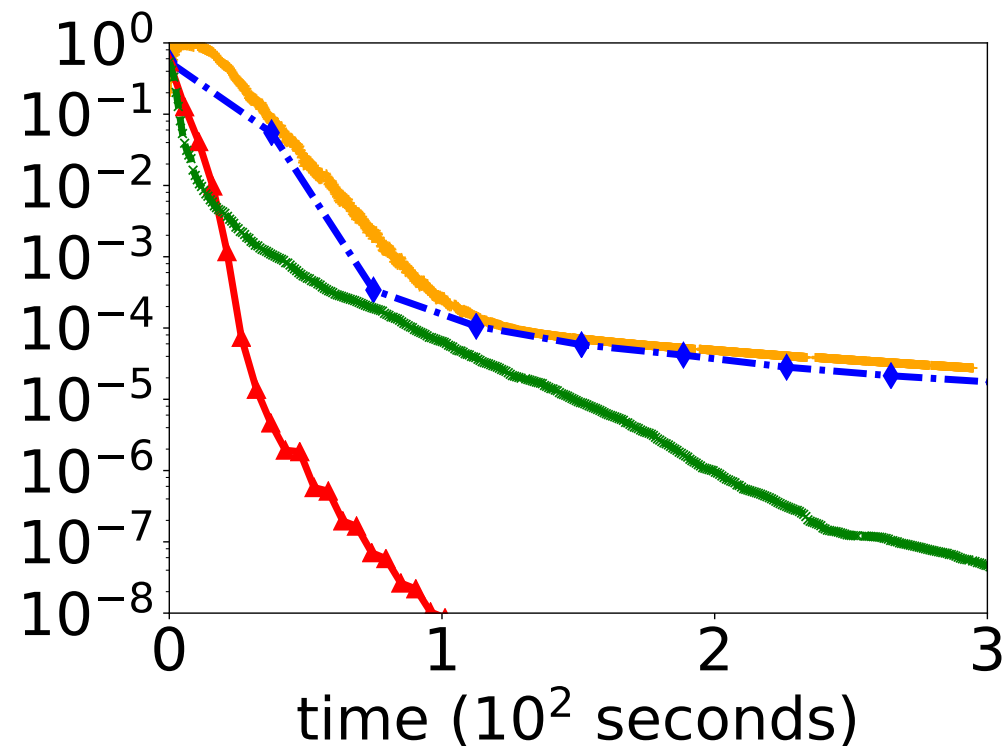ADMM · · · AGD · · · DANE — GIANT — · L-BFGS

AGD

# MNIST8M (n=1.6M, d=784)

- Digits "4" versus "9": 1.6M samples out of the total 8M samples

# MNIST8M (n=1.6M, d=784)

$\gamma = 10^{-4}$

$\gamma = 10^{-6}$

# How about Larger $d$?

- Split $\mathbf{X} \in \mathbb{R}^{n \times d}$ (by data) to $m = 59$ parts.

- Previously, local sample size $s = \dfrac{n}{m} \quad \gg \quad$ number of features $d$.

- Does GIANT work if $s \quad \approx \quad d$?

# Random Feature Maps (RFM)

- Generate <span style="color:red">10K</span> *random Fourier features* [Rahimi & Recht, 07] of the *RBF kernel*

$$k(\mathbf{x}_i, \mathbf{x}_j) \ = \ \exp\big( -\tfrac{1}{2\sigma} \|\mathbf{x}_i - \mathbf{x}_j\|_2^2 \big)$$
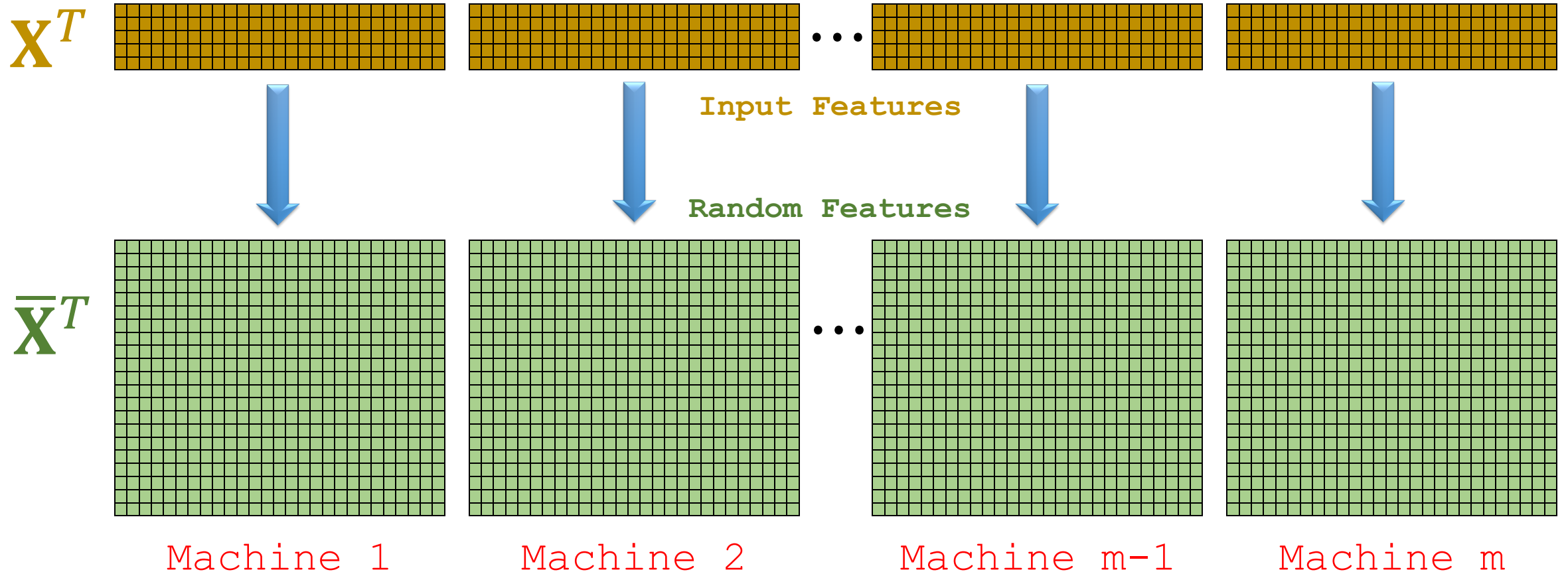
- Setting of RBF *kernel width parameter $\sigma$*:

$$\sigma \ = \ \frac{1}{n^2} \sum_{i=1}^{n} \sum_{j=1}^{n} \|\mathbf{x}_i - \mathbf{x}_j\|_2^2$$

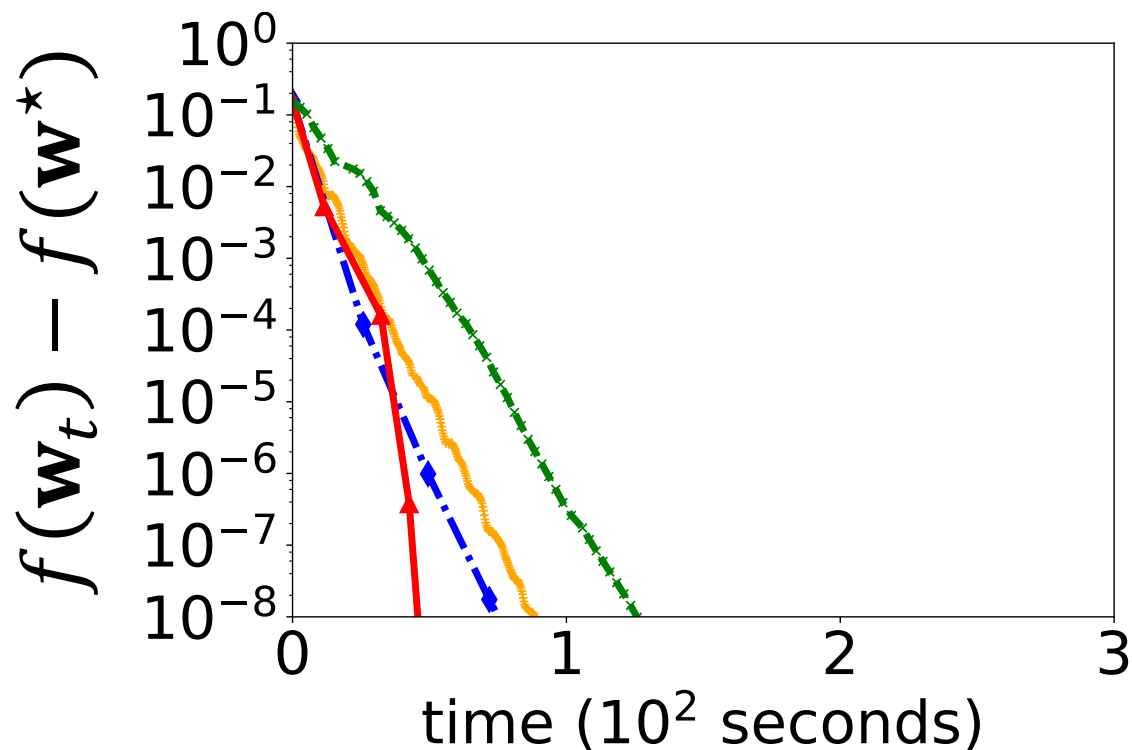- Replace the original features by the higher-dim random features.

**Reference**:
Rahimi & Recht. Random Features for Large-Scale Kernel Machines. In *NIPS,* 2007.
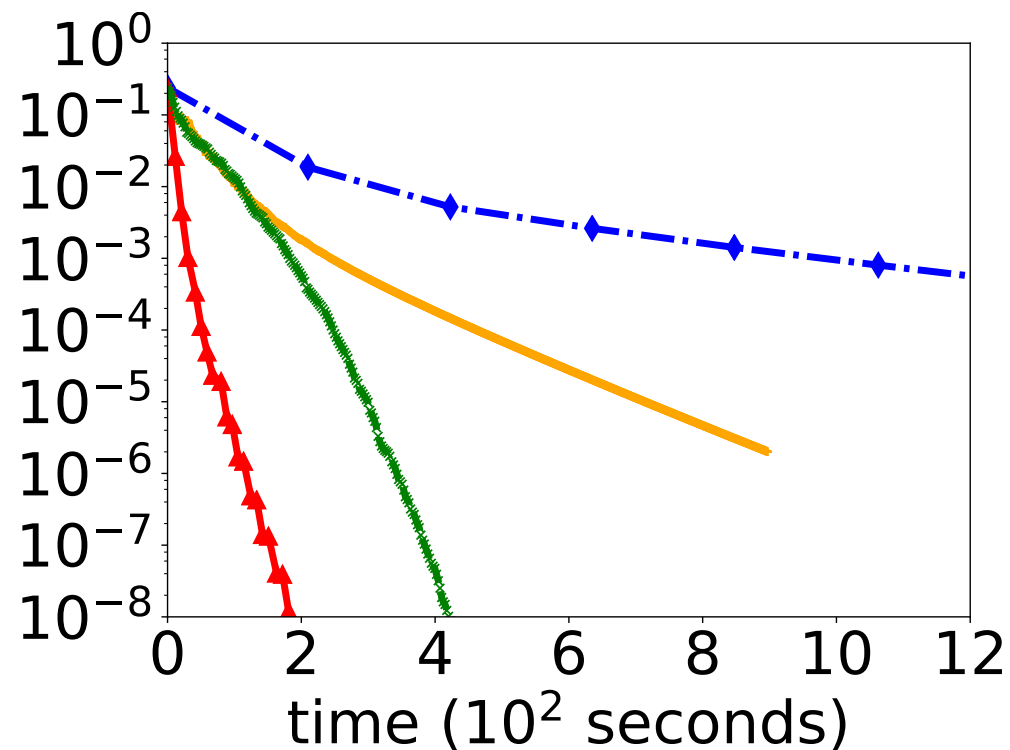
# Random Feature Maps (RFM)

Covtype with RFM (n=581K, d=10K)

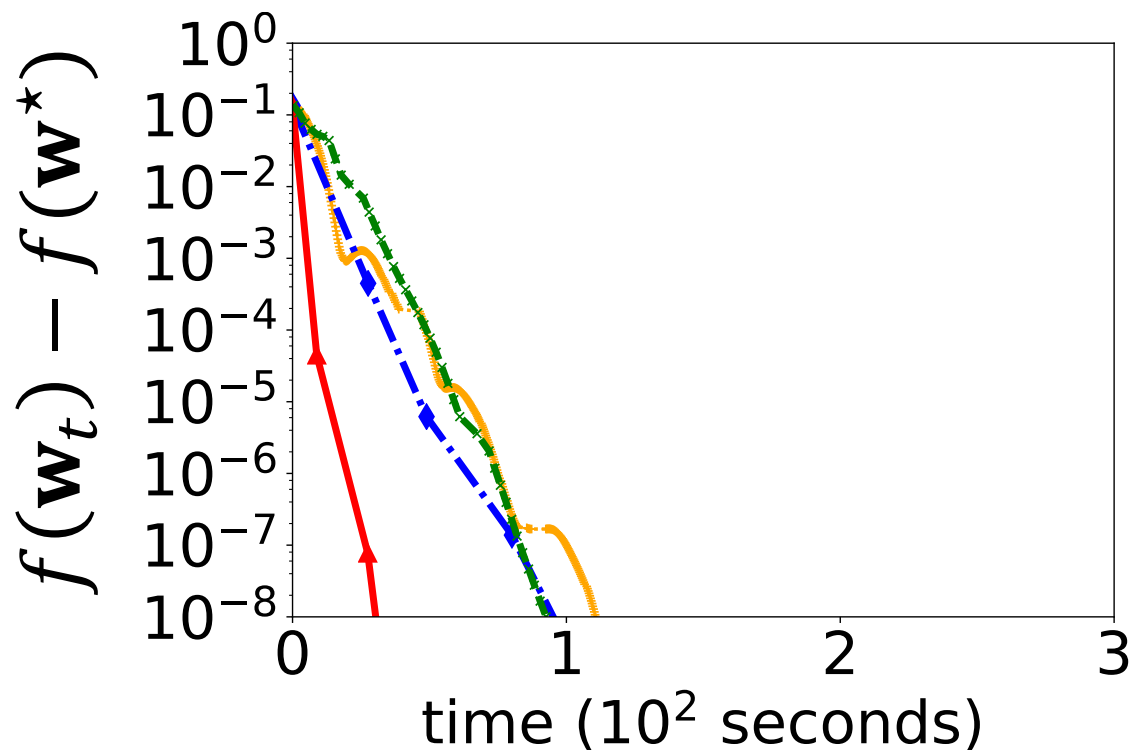$\gamma = 10^{-4}$
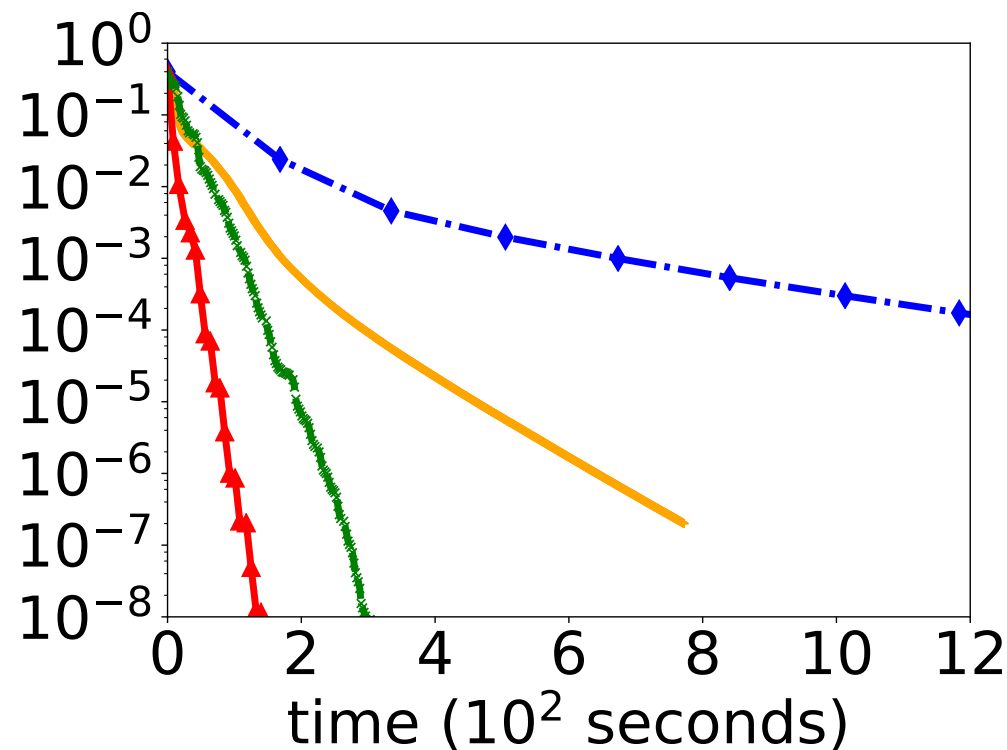
$\gamma = 10^{-6}$

ADMM   AGD   DANE   GIANT   L-BFGS

AGD

# Epsilon with RFM (n=400K, d=10K)

$\gamma = 10^{-4}$          $\gamma = 10^{-6}$

# MNIST8M with RFM (n=1.6M, d=10K)