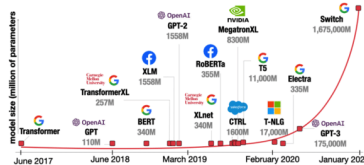


Aleksandr Beznosikov

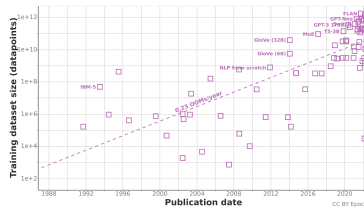
5 December 2023

# Contemporary learning challenges

- Exponential growth in model sizes and data volumes.



**Figure:** Dynamics of growth of modern language models



# Varieties of distributed learning

- Cluster learning (large players): we train within one large and powerful computing cluster
- Collaborative learning (all players): pooling computing resources over the Internet



## The most popular distributed setup

- Staging (horizontal, offline):

$$\min_{w \in \mathbb{R}^d} f(w) := \frac{1}{M} \sum_{m=1}^M f_m(w) := \frac{1}{M} \sum_{m=1}^M \frac{1}{n_m} \sum_{i=1}^{n_m} l(g(w, x_i), y_i).$$

- $w$  – model weights,  $g$  – model,  $l$  – loss function.
- The data is shared among  $M$  computing devices, each device  $m$  has its own local subsample  $\{x_i, y_i\}_{i=1}^{n_m}$  of size  $n_m$ .



# What are we fighting for?

- **Question:** distribution is necessary for parallelization, but why can't we achieve full parallelization?

# What are we fighting for?

- **Question:** distribution is necessary for parallelization, but why can't we achieve full parallelization?
- Communication costs are a waste of time.
- The communication bottleneck problem is relevant for all distributed learning productions.
- There are many ways to fight for effective communications.













## Unbiased compression: examples

- **Question:** What is  $\omega$  for random sparsification?

## Unbiased compression: examples

- Question:** What is  $\omega$  for random sparsification?  $\frac{d}{k}k$ .  
 Each coordinate takes part in  $\mathcal{Q}(x)$  with probability  $\frac{k}{d}$ , so

$$\begin{aligned}\mathbb{E} [\|\mathcal{Q}(x)\|^2] &= \mathbb{E} \left[ \sum_{i=1}^d [\mathcal{Q}(x)]_i^2 \right] \\ &= \frac{d^2}{k^2} \left[ \sum_{i=1}^d \frac{k}{d} [x]_i^2 \right] \\ &= \frac{d}{k} \|x\|^2.\end{aligned}$$

Here  $[\cdot]_i$  –  $i$ th coordinate of the vector.

## Three-level $\ell_2$ -quantization

$$\xi_i = \begin{cases} 1 & \text{with probability } \frac{|x_i|}{\|x\|_2}, \\ 0 & \text{with probability } 1 - \frac{|x_i|}{\|x\|_2}. \end{cases}$$





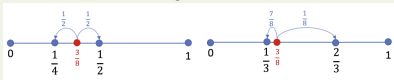

## Unbiased compression: examples

## Natural compression (random rounding to the power of two)

Consider the following operator:

$$[\mathcal{Q}(x)]_i = \begin{cases} \lfloor [x]_i \rfloor_2, & \text{with probability } p = \frac{[x]_i - \lceil [x]_i \rceil_2}{\lceil [x]_i \rceil_2 - \lfloor [x]_i \rfloor_2} \\ \lceil [x]_i \rceil_2, & \text{with probability } 1 - p \end{cases},$$

where  $[\cdot]_i$  –  $i$ -vector component,  $\lfloor \cdot \rfloor_2$  – is the nearest degree of two from the bottom,  $\lceil \cdot \rceil_2$  is the nearest degree of two from the top. We round to the two nearest powers of two, the probability of rounding is greater the closer the real number is to the corresponding power of two. It can be shown that  $\omega = \frac{9}{8}$ .



Horváth S. et al. Natural compression for distributed deep learning



## Quantized GD (QGD)

### Algorithm 1 QGD

**Input:** step size  $\gamma > 0$ , starting point  $w_0 \in \mathbb{R}^d$ , number of iterations  $K$

1: **for**  $k = 0, 1, \dots, K - 1$  **do**

2: Send  $w_k$  to all workers

- ▷ by server

3: **for**  $m = 1, \dots, M$  in parallel **do**4: Recieve  $w_k$  from server

- ▷ by workers

5: Compute  $\nabla f_m(w_k)$  in point  $w_k$

- ▷ by workers

6: Independently generate  $g_{k,m} = \mathcal{Q}(\nabla f_m(w_k))$

- ▷ by workers

7:       Send  $g_{k,m}$  to master

- ▷ by workers

8:     **end for**

9: Recieve  $g_{k,m}$  from all workers

- ▷ by server

10: Compute  $g_k = \frac{1}{M} \sum_{m=1}^M g_{k,m}$

- ▷ by server

11:  $w_{k+1} = w_k - \gamma g_k$

- ▷ by server

```
12: end for
```

**Output:**  $w^K$















## Unbiased compression: a proof

- We continue and apply the second property in the definition of compression:

$$\begin{aligned} \mathbb{E} [\|g_k\|^2 \mid w_k] &= \frac{1}{M^2} \sum_{m=1}^M \mathbb{E} [\|\mathcal{Q}(\nabla f_m(w_k))\|^2 \mid w_k] \\ &\quad + \frac{1}{M^2} \sum_{m \neq l} \langle \nabla f_m(w_k), \nabla f_l(w_k) \rangle \\ &\leq \frac{\omega}{M^2} \sum_{m=1}^M \|\nabla f_m(w_k)\|^2 + \|\nabla f(w_k)\|^2 \\ &\leq \frac{2\omega}{M^2} \sum_{m=1}^M \|\nabla f_m(w_k) - \nabla f_m(w^*)\|^2 \\ &\quad + \frac{2\omega}{M^2} \sum_{m=1}^M \|\nabla f_m(w^*)\|^2 + \|\nabla f(w_k) - \nabla f(w^*)\|^2. \end{aligned}$$













## QGD: convergence

## Theorem (QGD)

Let all local functions  $f_m$  be  $\mu$ -simply convex and have  $L$ -Lipschitz gradient, then if  $\eta \leq L^{-1} \left( \frac{2\omega}{M} + 1 \right)^{-1}$  then

$$\mathbb{E} [\|w_K - w^*\|^2] = \mathcal{O} \left( (1 - \gamma\mu)^K \|w_0 - w^*\|^2 + \frac{1}{K} \cdot \frac{2\omega}{\mu M^2} \sum_{m=1}^M \|\nabla f_m(w^*)\|^2 \right)$$

The selection of  $\gamma$  from the paper was also used to obtain this result:



Stich S. Unified Optimal Analysis of the (Stochastic) Gradient Method

- **Question:** what are the problems with this estimation? (recall the convergence estimate GD)

## QGD: convergence

### Theorem (QGD)

Let all local functions  $f_m$  be  $\mu$ -simply convex and have  $L$ -Lipschitz gradient, then if  $\eta \leq L^{-1} \left( \frac{2\omega}{M} + 1 \right)^{-1}$  then

$$\mathbb{E} [\|w_K - w^*\|^2] = \mathcal{O} \left( (1 - \gamma\mu)^K \|w_0 - w^*\|^2 + \frac{1}{K} \cdot \frac{2\omega}{\mu M^2} \sum_{m=1}^M \|\nabla f_m(w^*)\|^2 \right)$$

The selection of  $\gamma$  from the paper was also used to obtain this result:

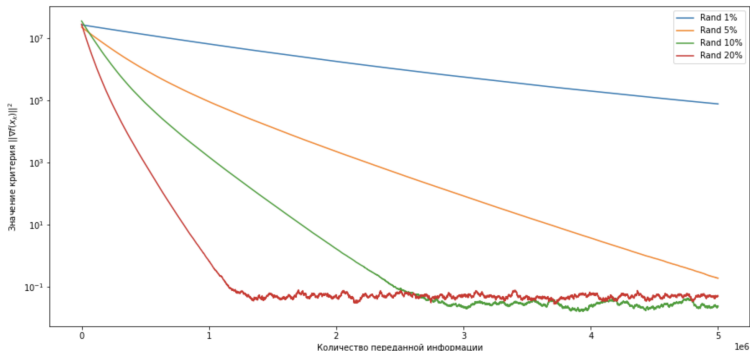


Stich S. Unified Optimal Analysis of the (Stochastic) Gradient Method

- **Question:** what are the problems with this estimation? (recall the convergence estimate GD) Sublinear convergence (depends on the heterogeneity of the data).

## QGD: convergence

- Behaviour in practice:



**Figure:** Behaviour of methods with unbiased compression operator and constant step size

- In theory the pitch was chosen cleverly, with constant step the theory predicts exactly the same effect – early plateauing.



## QGD and DIANA: convergence

- **Question:** What are some other questions about convergence/estimates of convergence?



## QGD and DIANA: convergence

- Compressors compress information  $\beta$  times and it is typical that  $\beta \geq \omega$ .

## QGD and DIANA: convergence

- Compressors compress information  $\beta$  times and it is typical that  $\beta \geq \omega$ .
- Best estimate on the number of information for the unaccelerated method with unbiased compression (DIANA):

$$\mathcal{O} \left( \left[ \frac{1}{\beta} + \frac{1}{M} \right] \frac{L}{\mu} \log \frac{1}{\varepsilon} \right).$$

- Estimate on the number of information for GD:

$$\mathcal{O} \left( \frac{L}{\mu} \log \frac{1}{\varepsilon} \right).$$



# QGD and DIANA: convergence

- Compressors compress information  $\beta$  times and it is typical that  $\beta \geq \omega$ .
- Best estimate on the number of information for the unaccelerated method with unbiased compression (DIANA):

$$\mathcal{O} \left( \left[ \frac{1}{\beta} + \frac{1}{M} \right] \frac{L}{\mu} \log \frac{1}{\varepsilon} \right).$$

- Estimate on the number of information for GD:

$$\mathcal{O} \left( \frac{L}{\mu} \log \frac{1}{\varepsilon} \right).$$

- The unbiased compressor provably improves the number of transmitted information, an improvement factor:  $\left[ \frac{1}{\beta} + \frac{1}{M} \right]$ .



# Centralised communications without a server

Operation	Before				After			
Broadcast	Node 0 $x$	Node 1	Node 2	Node 3	Node 0 $x$	Node 1 $x$	Node 2 $x$	Node 3 $x$
Reduce(-to-one)	Node 0 $x^{(0)}$	Node 1 $x^{(1)}$	Node 2 $x^{(2)}$	Node 3 $x^{(3)}$	Node 0 $\sum_j x^{(j)}$	Node 1	Node 2	Node 3
Scatter	Node 0 $x_0$ $x_1$ $x_2$ $x_3$	Node 1	Node 2	Node 3	Node 0 $x_0$	Node 1 $x_1$	Node 2 $x_2$	Node 3 $x_3$
Gather	Node 0 $x_0$	Node 1 $x_1$	Node 2 $x_2$	Node 3 $x_3$	Node 0 $x_0$ $x_1$ $x_2$ $x_3$	Node 1	Node 2	Node 3
Allgather	Node 0 $x_0$	Node 1 $x_1$	Node 2 $x_2$	Node 3 $x_3$	Node 0 $x_0$ $x_1$ $x_2$ $x_3$	Node 1 $x_0$ $x_1$ $x_2$ $x_3$	Node 2 $x_0$ $x_1$ $x_2$ $x_3$	Node 3 $x_0$ $x_1$ $x_2$ $x_3$
Reduce-scatter	Node 0 $x_0^{(0)}$ $x_1^{(0)}$ $x_2^{(0)}$ $x_3^{(0)}$	Node 1 $x_0^{(1)}$ $x_1^{(1)}$ $x_2^{(1)}$ $x_3^{(1)}$	Node 2 $x_0^{(2)}$ $x_1^{(2)}$ $x_2^{(2)}$ $x_3^{(2)}$	Node 3 $x_0^{(3)}$ $x_1^{(3)}$ $x_2^{(3)}$ $x_3^{(3)}$	Node 0 $\sum_j x_0^{(j)}$	Node 1 $\sum_j x_1^{(j)}$	Node 2 $\sum_j x_2^{(j)}$	Node 3 $\sum_j x_3^{(j)}$
Allreduce	Node 0 $x^{(0)}$	Node 1 $x^{(1)}$	Node 2 $x^{(2)}$	Node 3 $x^{(3)}$	Node 0 $\sum_j x^{(j)}$	Node 1 $\sum_j x^{(j)}$	Node 2 $\sum_j x^{(j)}$	Node 3 $\sum_j x^{(j)}$

# Ring AllReduce

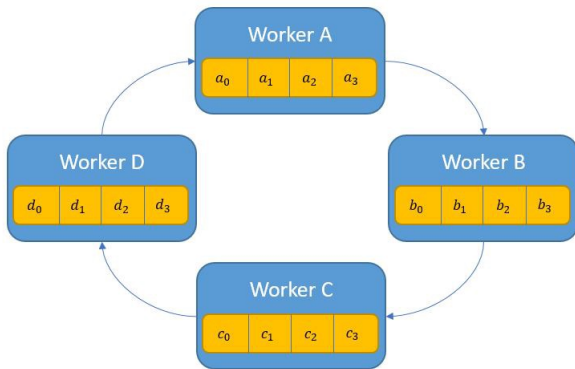


Figure: Picture from here

# Ring AllReduce: first step of averaging

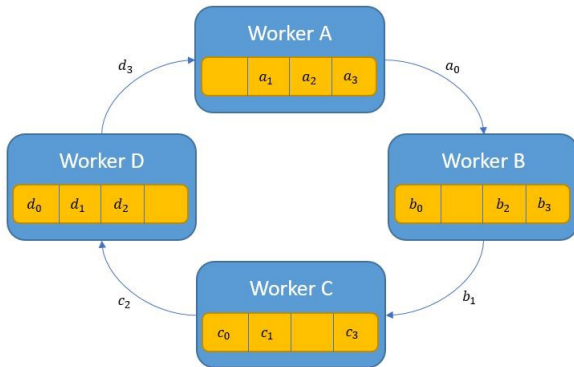


Figure: Picture from here





# Ring AllReduce: final

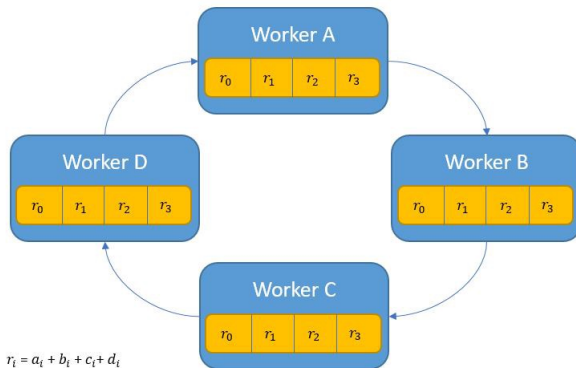


Figure: Картинка отсюда





### Algorithm 1 QGD

1: **for**  $k = 0, 1, \dots, K - 1$  **do**2: **for**  $m = 1, \dots, M$  in parallel **do**

3: Compute  $\nabla f_m(w_k)$  in  $w_k$

4: Independently generate  $g_{k,m} = \mathcal{Q}(\nabla f_m(w_k))$

5: Run AllReduce  $\{g_{k,m}\}$  and get  $g_k = \frac{1}{M} \sum_{m=1}^M g_{k,m}$

6:  $w_{k+1} = w_k - \gamma g_k$

```
7:   end for
```

8: end for

**Output:**  $w^K$

# Quantized GD (QGD) with AllReduce

---

## Algorithm 1 QGD

---

**Input:** step size  $\gamma > 0$ , starting point  $w_0 \in \mathbb{R}^d$ , number of iterations  $K$

```
1: for  $k = 0, 1, \dots, K - 1$  do
2:   for  $m = 1, \dots, M$  in parallel do
3:     Compute  $\nabla f_m(w_k)$  in  $w_k$ 
4:     Independently generate  $g_{k,m} = \mathcal{Q}(\nabla f_m(w_k))$ 
5:     Run AllReduce  $\{g_{k,m}\}$  and get  $g_k = \frac{1}{M} \sum_{m=1}^M g_{k,m}$ 
6:      $w_{k+1} = w_k - \gamma g_k$ 
7:   end for
8: end for
```

**Output:**  $w^K$

---

- **Question:** What kind of problems might appear with (for example) Rank? The same non-zero coordinates on different devices can cause collisions.











# Offset compression: examples

- Various examples of compressors (sparsifiers, rounding, etc):
  -  Beznosikov A. et al. On Biased Compression for Distributed Learning
- A practical biased compressor based on iterative SVD decomposition:
  -  Vogels T. et al. PowerSGD: Practical Low-Rank Gradient Compression for Distributed Optimization



# Biased compression: idea and proof in the case of 1 node

- Use the same approach as in the unbiased case (QGD).



# Biased compression: idea and proof in the case of 1 node

- Use the same approach as in the unbiased case (QGD).
- Let's prove it in the case of a single node:

$$w_{k+1} = w_k - \gamma C(\nabla f(w_k)).$$

Let  $f$  have  $L$ -Lipschitz gradient and be  $\mu$ -simply convex.

- Let's start by taking advantage of the Lipschitzness of the gradient:

$$\begin{aligned} f(w_{k+1}) &= f(w_k - \gamma C(\nabla f(w_k))) \\ &\leq f(w_k) + \langle \nabla f(w_k), -\gamma C(\nabla f(w_k)) \rangle + \frac{L}{2} \| -\gamma C(\nabla f(w_k)) \|^2 \\ &= f(w_k) - \gamma \langle C(\nabla f(w_k)), \nabla f(w_k) \rangle + \frac{\gamma^2 L}{2} \| C(\nabla f(w_k)) \|^2. \end{aligned}$$

- Compressor definition:

$$\begin{aligned} & \|\nabla f(w_k)\|^2 - 2\mathbb{E}_C [\langle C(\nabla f(w_k)), \nabla f(w_k) \rangle] + \mathbb{E}_C [\|C(\nabla f(w_k))\|^2] \\ &= \mathbb{E}_C [\|C(\nabla f(w_k)) - \nabla f(w_k)\|^2] \leq \left(1 - \frac{1}{\delta}\right) \|\nabla f(w_k)\|^2. \end{aligned}$$

# Biased compression: idea and proof in the case of 1 node

- Compressor definition:

$$\begin{aligned} \|\nabla f(w_k)\|^2 - 2\mathbb{E}_C [\langle C(\nabla f(w_k)), \nabla f(w_k) \rangle] + \mathbb{E}_C [\|C(\nabla f(w_k))\|^2] \\ = \mathbb{E}_C [\|C(\nabla f(w_k)) - \nabla f(w_k)\|^2] \leq \left(1 - \frac{1}{\delta}\right) \|\nabla f(w_k)\|^2. \end{aligned}$$

- From where:

$$-\gamma \mathbb{E}_C [\langle C(\nabla f(w_k)), \nabla f(w_k) \rangle] + \frac{\gamma}{2} \mathbb{E}_C [\|C(\nabla f(w_k))\|^2] \leq -\frac{\gamma}{2\delta} \|\nabla f(w_k)\|^2.$$



# Biased compression: idea and proof in the case of 1 node

- From the previous two slides:

$$f(w_{k+1}) - \leq f(w_k) - \gamma \langle C(\nabla f(w_k)), \nabla f(w_k) \rangle + \frac{\gamma^2 L}{2} \|C(\nabla f(w_k))\|^2.$$

$$- \gamma \mathbb{E}_C [\langle C(\nabla f(w_k)), \nabla f(w_k) \rangle] + \frac{\gamma}{2} \mathbb{E}_C [\|C(\nabla f(w_k))\|^2] \leq - \frac{\gamma}{2\delta} \|\nabla f(w_k)\|^2.$$

- Add, subtract  $f(w^*)$  from both parts, and take the full expectation matrix:

$$\begin{aligned} \mathbb{E} [f(w_{k+1}) - f(w^*)] &\leq \mathbb{E} [f(w_k) - f(w^*)] - \frac{\gamma}{2} (1 - \gamma L) \mathbb{E} [\|C(\nabla f(w_k))\|^2] \\ &\quad - \frac{\gamma}{2\delta} \mathbb{E} [\|\nabla f(w_k)\|^2]. \end{aligned}$$





- From the previous slide:

$$\mathbb{E}[f(w_{k+1}) - f(w^*)] \leq \mathbb{E}[f(w_k) - f(w^*)] - \frac{\gamma}{2\delta} \mathbb{E}[\|\nabla f(w_k)\|^2].$$

# Biased compression: idea and proof in the case of 1 node

- From the previous slide:

$$\mathbb{E}[f(w_{k+1}) - f(w^*)] \leq \mathbb{E}[f(w_k) - f(w^*)] - \frac{\gamma}{2\delta} \mathbb{E}[\|\nabla f(w_k)\|^2].$$

- Strong convexity (or even the weaker PL condition):

$$2\mu(f(w_k) - f(w^*)) \leq \|\nabla f(w_k)\|^2.$$



# Biased compression: a theorem in the case of 1 node

Theorem (convergence of QGD with shifted compression in case of 1 node)

Let  $f$   $\mu$  be strongly convex (or PL) and have  $L$ -Lipschitz gradient, then QGD for one node with step  $\gamma \leq 1/L$  and with biased compressor with parameter  $\delta$  converges and is satisfied:

$$f(w_K) - f(w^*) \leq \left(1 - \frac{\gamma\mu}{\delta}\right)^K (f(w_0) - f(w^*)).$$



Beznosikov A. et al. On Biased Compression for Distributed Learning











## Biased compression: it's not that simple

- Consider the following distributed problem with  $M = 3$ ,  $d = 3$  and local functions:

$$f_1(w) = \langle a, w \rangle^2 + \frac{1}{4} \|w\|^2, \quad f_2(w) = \langle b, w \rangle^2 + \frac{1}{4} \|w\|^2, \quad f_3(w) = \langle c, w \rangle^2 + \frac{1}{4} \|w\|^2$$

where  $a = (-3, 2, 2)$ ,  $b = (2, -3, 2)$  и  $c = (2, 2, -3)$ .

- Question:** where is her optimum?  $(0, 0, 0)$ .
- Let the starting point  $w_0 = (t, t, t)$  for some  $t > 0$ . Then the local gradients are:

$$\nabla f_1(w_0) = \frac{t}{2}(-11, 9, 9), \quad \nabla f_2(w_0) = \frac{t}{2}(9, -11, 9), \quad \nabla f_3(w_0) = \frac{t}{2}(9, 9, -11).$$

- Question:** what will the QGD (gradient descent with compressions) step look like if we use *Top1* compression?

$$w_1 = (t, t, t) + \eta \cdot \frac{11}{6} (t, t, t) = \left(1 + \frac{11\eta}{6}\right) w_0.$$

- We move away from the solution geometrically for any  $\eta > 0$ .





## Offset compression: error compensation

- Let's try to remember what we didn't pass on in the communication process:

$$e_{1,m} = 0 + \gamma \nabla f_m(w_0) - C(0 + \gamma \nabla f_m(w_0)).$$

- And add this to future parcels:

$$C(e_{1,m} + \gamma \nabla f_m(w_1))$$

- In an arbitrary iteration, it is written as follows:

Parcel:  $C(e_{k,m} + \gamma \nabla f_m(w_k)),$

$$e_{k+1,m} = e_{k,m} + \gamma \nabla f_m(w_k) - C(e_{k,m} + \gamma \nabla f_m(w_k))$$

# Offset compression: error compensation

- Let's try to remember what we didn't pass on in the communication process:

$$e_{1,m} = 0 + \gamma \nabla f_m(w_0) - C(0 + \gamma \nabla f_m(w_0)).$$

- And add this to future parcels:

$$C(e_{1,m} + \gamma \nabla f_m(w_1))$$

- In an arbitrary iteration, it is written as follows:

$$\text{Parcel: } C(e_{k,m} + \gamma \nabla f_m(w_k)),$$

$$e_{k+1,m} = e_{k,m} + \gamma \nabla f_m(w_k) - C(e_{k,m} + \gamma \nabla f_m(w_k))$$

- This technique is called error compensation (error feedback).



Stich S. et al. Sparsified SGD with memory

## QGD with error feedback

---

**Algorithm 1** QGD c error feedback

**Input:** step size  $\gamma > 0$ , starting point  $w_0 \in \mathbb{R}^d$ , starting errors  $e_{0,m} = 0$  for all  $m$  from 1 to  $M$ , number of iterations  $K$

```

1: for  $k = 0, 1, \dots, K - 1$  do
2:   Send  $x_k$  to all workers ▷ by server
3:   for  $m = 1, \dots, M$  in parallel do
4:     Recieve  $w_k$  from the master ▷ by workers
5:     Compute  $\nabla f(w_k)$  in  $w_k$  ▷ by workers
6:     Generate  $g_{k,m} = C(e_{k,m} + \gamma \nabla f(w_k))$  ▷ by workers
7:     Вычислить  $e_{k+1,m} = e_{k,m} + \gamma \nabla f_m(w_k) - g_{k,m}$  ▷ by workers
8:     Send  $g_{k,m}$  to master ▷ by workers
9:   end for
10:  Recieve  $g_{k,m}$  from all workers ▷ by server
11:  Compute  $g_k = \frac{1}{M} \sum_{m=1}^M g_{k,m}$  ▷ by server
12:   $w_{k+1} = w_k - g_k$  ▷ by server
13: end for

```

**Output:**  $w_{12}$

# QGD with error feedback: convergence

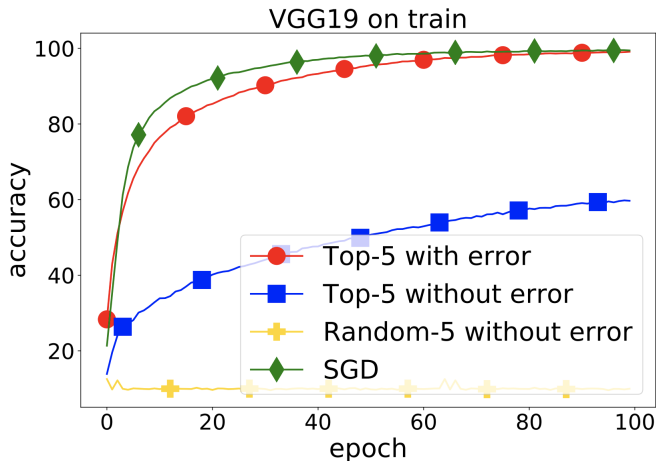


Figure: Accuracy during training of VGG19 on CIFAR10 using different compressors





# Biased compression: solving the plateau issue

- Idea like DIANA: memory + compression of difference

---

## Algorithm 1 EF21 (sketch)

---

- Each device  $m$  possesses a vector of "memory"  $g_0^m = 0$
  - The server stores  $h_0 = \frac{1}{M} \sum_{m=1}^M h_0^m = 0$
  - Send a compressed version of the difference to the server  $C(\nabla f_m(w^k) - h_k^m)$
  - Refreshing the memory  $h_{k+1}^m = h_k^m + C(\nabla f_m(w^k) - h_k^m)$
  - Server compute  $h_{k+1} = h_k + \frac{1}{M} \sum_{m=1}^M C(\nabla f_m(w^k) - h_k^m)$
  - Update:  $w_{k+1} = w_k - \gamma h_{k+1}$
- 



Richtarik P. et al. EF21: A New, Simpler, Theoretically Better, and Practically Faster Error Feedback

# Unbiased vs biased

- Best estimate on the number of communications for the unaccelerated method with unbiased compression (DIANA):

$$\mathcal{O}\left(\left[1 + \frac{\omega}{M}\right] \frac{L}{\mu} \log \frac{1}{\varepsilon}\right).$$

- Best estimate on the number of communications for the unaccelerated method with biased compression (EF-21):

$$\mathcal{O}\left([1+\delta]\frac{L}{\mu}\log\frac{1}{\varepsilon}\right).$$

- It has already been discussed that these estimates are worse than for the baseline GD.

# Unbiased vs biased

- Compressors compress information  $\beta$  times and it is typical that  $\beta \geq \omega$  and  $\beta \geq \delta$ .



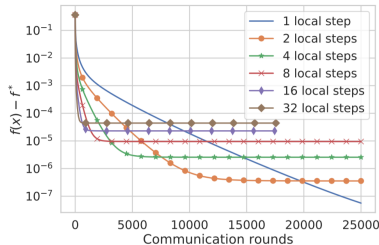







# Convergence

- Typical convergence of this type of methods:



**Figure:** Convergence of the Local Method in practice for logistic regression.

- Faster in terms of communications, worse quality of marginal accuracy.
-  Khaled A. et al. Tighter Theory for Local SGD on Identical and Heterogeneous Data













## What do we want to achieve, anyway?

- Lower bounds:

$$K = \Omega \left( \sqrt{\frac{L}{\mu}} \log \frac{1}{\varepsilon} \right).$$

 $L$  и  $\mu$  – smoothness and strong convexity constants  $f$ .

- **Question:** which method will give these estimates? Distributed version of Nesterov's method with 1 local step between communications.
- We note that local steel methods for stochastic productions.

## What do we want to achieve, anyway?

- Lower bounds:

$$K = \Omega \left( \sqrt{\frac{L}{\mu}} \log \frac{1}{\varepsilon} \right).$$

 $L$  и  $\mu$  – smoothness and strong convexity constants  $f$ .

- **Question:** which method will give these estimates? Distributed version of Nesterov's method with 1 local step between communications.
- We note that local steel methods for stochastic productions.
- But even here, there is generally no improvement.  
Woodworth B. The Min-Max Complexity of Distributed Stochastic Convex Optimization with Intermittent Communication
- But there are productions where localised methods shoot out.



## Data similarity







# Data similarity

- A distributed learning task:

$$f(w) = \frac{1}{M} \sum_{m=1}^M f_m(w) = \frac{1}{M} \sum_{m=1}^M \left[ \frac{1}{N} \sum_{i=1}^N \ell(w, z_i) \right],$$

where  $z_i$  is the sampling element  $(x_i, y_i)$ ,  $\ell$  is the loss function (it has  $f$  and  $g$  sewn into it).

- Suppose we can split the training sample evenly across devices (e.g., if cluster or collaborative computing on open data is used).
- This gives the similarity of the local loss functions.

# Data similarity

- A distributed learning task:

$$f(w) = \frac{1}{M} \sum_{m=1}^M f_m(w) = \frac{1}{M} \sum_{m=1}^M \left[ \frac{1}{N} \sum_{i=1}^N \ell(w, z_i) \right],$$

where  $z_i$  is the sampling element  $(x_i, y_i)$ ,  $\ell$  is the loss function (it has  $f$  and  $g$  sewn into it).

- Suppose we can split the training sample evenly across devices (e.g., if cluster or collaborative computing on open data is used).
- This gives the similarity of the local loss functions.
- It is argued that for any  $w$

$$\|\nabla^2 f_m(w) - \nabla^2 f(w)\| \leq \delta.$$







## Similarity parameter: bottom line

- As a result, we have

$$\|\nabla^2 f_m(w) - \nabla^2 f(w)\| \leq \delta \sim \frac{L}{\sqrt{N}}.$$

- For quadratic problems, we can obtain an estimate of the form:

$$\|\nabla^2 f_m(w) - \nabla^2 f(w)\| \leq \delta \sim \frac{L}{N}.$$



Hendrikx H. et al. Statistically Preconditioned Accelerated Gradient Method for Distributed Optimization

## Similarity parameter: bottom line

- As a result, we have

$$\|\nabla^2 f_m(w) - \nabla^2 f(w)\| \leq \delta \sim \frac{L}{\sqrt{N}}.$$

- For quadratic problems, we can obtain an estimate of the form:

$$\|\nabla^2 f_m(w) - \nabla^2 f(w)\| \leq \delta \sim \frac{L}{N}.$$



Hendriks H. et al. Statistically Preconditioned Accelerated Gradient Method for Distributed Optimization

- In any case, the conclusion follows: the larger the local sample size, the smaller the similarity parameter (hessians are similar to each other).





## Method in general terms

- Consider the mirror descent:

$$w_{k+1} = \arg \min_{w \in \mathbb{R}^d} (\gamma \langle \nabla f(w_k), w \rangle + V(w, w_k)),$$

where  $V(x, y)$  is the Bregman divergence generated by a strictly-convex function  $\varphi(x)$ :

$$V(x, y) = \varphi(x) - \varphi(y) - \langle \nabla \varphi(y); x - y \rangle.$$

- **Question:** Which method will result if  $\varphi(x) = \frac{1}{2}\|x\|^2$ ?







## Convergence in general: a proof

- The first optimality condition for the mirror descent step:

$$\gamma \nabla f(w_k) + \nabla \varphi(w_{k+1}) - \nabla \varphi(w_k) = 0.$$

- From it (here  $w^*$  – optimum):

$$\langle \gamma \nabla f(w_k) + \nabla \varphi(w_{k+1}) - \nabla \varphi(w_k), w_{k+1} - w^* \rangle = 0.$$

$$\begin{aligned} \langle \gamma \nabla f(w_k), w^{k+1} - w^* \rangle &= \langle \nabla \varphi(w_k) - \nabla \varphi(w_{k+1}), w^{k+1} - w^* \rangle \\ &= V(w^*, w_k) - V(w^*, w_{k+1}) - V(w_{k+1}, w_k). \end{aligned}$$

(the last statement is called Pythagoras' theorem for Bregman divergences and is verified by definition)







## Convergence in general: a proof

- Подставим  $\gamma = \frac{1}{L_\varphi}$ :

$$\begin{aligned} \langle \nabla f(w_k), w^{k+1} - w^k \rangle + L_\varphi V(w_{k+1}, w_k) \\ = L_\varphi V(w^*, w_k) - L_\varphi V(w^*, w_{k+1}) \\ - \langle \nabla f(w_k), w_k - w^* \rangle. \end{aligned}$$

- Let us use the definition of smoothness with respect to  $\varphi \in x = w_{k+1}$ ,  $y = w_k$ :

$$f(w_{k+1}) - f(w_k) \leq \langle \nabla f(w_k); w_{k+1} - w_k \rangle + L_\varphi V(w_{k+1}, w_k).$$

## Convergence in general: a proof

- Подставим  $\gamma = \frac{1}{L_\varphi}$ :

$$\begin{aligned} & \langle \nabla f(w_k), w^{k+1} - w^k \rangle + L_\varphi V(w_{k+1}, w_k) \\ &= L_\varphi V(w^*, w_k) - L_\varphi V(w^*, w_{k+1}) \\ & \quad - \langle \nabla f(w_k), w_k - w^* \rangle. \end{aligned}$$

- Let us use the definition of smoothness with respect to  $\varphi \in X = W_{k+1}$ ,  $Y = W_k$ :

$$f(w_{k+1}) - f(w_k) \leq \langle \nabla f(w_k); w_{k+1} - w_k \rangle + L_\varphi V(w_{k+1}, w_k).$$

- Let's combine the previous two:

$$f(w_{k+1}) - f(w_k) \leq L_\varphi V(w^*, w_k) - L_\varphi V(w^*, w_{k+1}) - \langle \nabla f(w_k), w_k - w^* \rangle$$









## Convergence in general: theorem

## Theorem (convergence of mirror descent)

Let  $\varphi$  and  $f$  satisfy the definition above, then the mirror descent with step  $\gamma = \frac{1}{L_\varphi}$  converges and is satisfied:

$$V(w^*, w_K) \leq \left(1 - \frac{\mu_\varphi}{L_\varphi}\right)^K V(w^*, w_0).$$



Lu H. et al. Relatively-Smooth Convex Optimization by First-Order Methods, and Applications





## Method for data similarity

- Mirror descent:

$$w_{k+1} = \arg \min_{w \in \mathbb{R}^d} (\gamma \langle \nabla f(w_k), w \rangle + V(w, w_k)),$$

where  $V(x, y)$  is the Bregman divergence generated by the function  $\varphi(x)$  (here we need to require that  $f_1$  is convex):

$$\varphi(x) = f_1(x) + \frac{\delta}{2} \|x\|^2.$$

The function  $f_1$  is stored on the server.

- **Question:** What is the number of communications that occur in  $K$  iterations of such a mirror descent?

## Method for data similarity

- Mirror descent:

$$w_{k+1} = \arg \min_{w \in \mathbb{R}^d} (\gamma \langle \nabla f(w_k), w \rangle + V(w, w_k)),$$

where  $V(x, y)$  is the Bregman divergence generated by the function  $\varphi(x)$  (here we need to require that  $f_1$  is convex):

$$\varphi(x) = f_1(x) + \frac{\delta}{2} \|x\|^2.$$

The function  $f_1$  is stored on the server.

- Question:** What is the number of communications that occur in  $K$  iterations of such a mirror descent?  $K$  of communications (the number of gradient counts  $\nabla f$ ), computing  $\arg \min$  requires only computations on the server.



## Convergence for data similarity: a proof

- Recall that convergence is defined in terms of constants from the relation:

$$\mu_\varphi \nabla^2 \varphi(w) \preceq \nabla^2 f(w) \preceq L_\varphi \nabla^2 \varphi(w),$$





















## Theorem (convergence for data similarity)

Let  $f$  be strongly convex,  $f_i$  be convex, and  $\ell$  be smooth, and  $\varphi(w) = f_1(w) + \delta \|w\|^2$ , then mirror descent with step  $\gamma = 1$  converges and is satisfied:

$$V(w^*, w_K) \leq \left(1 - \frac{\mu}{\mu + 2\delta}\right)^K V(w^*, w_0).$$

## Convergence for data similarity: theorem

## Theorem (convergence for data similarity)

Let  $f$  be strongly convex,  $f_i$  be convex, and  $\ell$  be smooth, and  $\varphi(w) = f_1(w) + \delta \|w\|^2$ , then mirror descent with step  $\gamma = 1$  converges and is satisfied:

$$V(w^*, w_K) \leq \left(1 - \frac{\mu}{\mu + 2\delta}\right)^K V(w^*, w_0).$$

- It means that if we need to achieve an accuracy  $\varepsilon$  ( $V(w^*, w_K) \sim \varepsilon$ ), then we need to

$$K = \left( \left[ 1 + \frac{\delta}{\mu} \right] \log \frac{V(w^*, w_0)}{\varepsilon} \right) \text{ communications.}$$

# Better?

- The estimate on the number of communications under data similarity:

$$K = \mathcal{O} \left( \left[ 1 + \frac{\delta}{\mu} \right] \log \frac{1}{\varepsilon} \right).$$

- The estimate on the number of communications for the basic distributed gradient descent:

$$K = \mathcal{O} \left( \frac{L}{\mu} \log \frac{1}{\varepsilon} \right).$$





## Another look at the mirror escapement

- Mirror descent with  $\gamma = 1$ :

$$w_{k+1} = \arg \min_{w \in \mathbb{R}^d} (\langle \nabla f(w_k), w \rangle + V(w, w_k)),$$

where  $V(x, y)$  is the Bregman divergence generated by the function  $\varphi(x)$ :

$$\varphi(x) = f_1(x) + \frac{\delta}{2}\|x\|^2.$$

- Substitute  $\varphi(x)$ :

$$w_{k+1} = \arg \min_{w \in \mathbb{R}^d} \left( f_1(w) + \langle \nabla f(w_k) - \nabla f_1(w_k), w \rangle + \frac{\delta}{2} \|w - w_k\|^2 \right).$$

- Or:

$$w_{k+1} = \arg \min_{w \in \mathbb{R}^d} \left( \frac{1}{\delta} f_1(w) + \frac{1}{2} \left\| w - \left( w_k - \frac{1}{\delta} (\nabla f(w_k) - \nabla f_1(w_k)) \right) \right\|^2 \right)$$

Bottom line about mirror descent

- The idea of regularizing the local subproblem came up.
- The idea of sliding  $\approx$  proximal method with inexactness.
- Proximal method for a composite target function  $g_1(w) + g_2(w)$ :

$$w_{k+1} = \arg \min_{w \in \mathbb{R}^d} \left( \gamma g_2(w) + \frac{1}{2} \|w - (w_k - \gamma g_1(w_k))\|^2 \right).$$

- In our case,  $g_1 = f - f_1$ ,  $g_2 = f_1$ .



- This problem has quite a history:

Reference	Communication complexity	Local gradient complexity	Order	Limitations
DANE [42]	$\mathcal{O}\left(\frac{\delta^2}{\mu^2} \log \frac{1}{\epsilon}\right)$	$—^{(2)}$	1st	quadratic
DiSCO [51]	$\mathcal{O}\left(\sqrt{\frac{\delta}{\mu}} \left(\log \frac{1}{\epsilon} + C^2 \Delta F_0\right) \log \frac{L}{\mu}\right)$	$\mathcal{O}\left(\sqrt{\frac{\delta}{\mu}} \left(\log \frac{1}{\epsilon} + C^2 \Delta F_0\right) \log \frac{L}{\mu}\right)$	2nd	$C$ -self-concordant <sup>(3)</sup>
AIDE [40]	$\mathcal{O}\left(\sqrt{\frac{\delta}{\mu}} \log \frac{1}{\epsilon} \log \frac{L}{\delta}\right)$	$\mathcal{O}\left(\sqrt{\frac{\delta}{\mu}} \sqrt{\frac{L}{\mu}} \log \frac{1}{\epsilon} \log \frac{L}{\delta}\right)^{(4)}$	1st	quadratic
DANE-LS [50]	$\mathcal{O}\left(\frac{\delta}{\mu} \log \frac{1}{\epsilon}\right)$	$\mathcal{O}\left(\sqrt{\frac{\delta}{\mu}} \frac{3^{3/2}}{3/2} \log \frac{1}{\epsilon}\right)^{(5)}$	1st/2nd	quadratic <sup>(6)</sup>
DANE-HB [50]	$\mathcal{O}\left(\sqrt{\frac{\delta}{\mu}} \log \frac{1}{\epsilon}\right)$	$\mathcal{O}\left(\sqrt{\frac{\delta}{\mu}} \frac{\delta}{\mu} \log \frac{1}{\epsilon}\right)^{(5)}$	1st/2nd	quadratic <sup>(6)</sup>
SONATA [45]	$\mathcal{O}\left(\frac{\delta}{\mu} \log \frac{1}{\epsilon}\right)$	$—^{(2)}$	1st	decentralized
SPAG [21]	$\mathcal{O}\left(\sqrt{\frac{\delta}{\mu}} \log \frac{1}{\epsilon}\right)^{(1)}$	$—^{(2)}$	1st	$M$ -Lipshitz hessian
DiRegINA [12]	$\mathcal{O}\left(\frac{\delta}{\mu} \log \frac{1}{\epsilon} + \sqrt{\frac{M\delta R_0}{\mu}}\right)$	$—^{(2)}$	2nd	$M$ -Lipshitz hessian
ACN [1]	$\mathcal{O}\left(\sqrt{\frac{\delta}{\mu}} \log \frac{1}{\epsilon} + \sqrt[3]{\frac{M\delta R_0}{\mu}}\right)$	$—^{(2)}$	2nd	$M$ -Lipshitz hessian
Acc SONATA [46]	$\mathcal{O}\left(\sqrt{\frac{\delta}{\mu}} \log \frac{1}{\epsilon} \log \frac{\delta}{\mu}\right)$	$—^{(2)}$	1st	decentralized
This paper	$\mathcal{O}\left(\sqrt{\frac{\delta}{\mu}} \log \frac{1}{\epsilon}\right)$	$\mathcal{O}\left(\sqrt{\frac{L}{\mu}} \log \frac{1}{\epsilon}\right)$	1st	

In particular, the mirror descent approach with unusual divergence is called DANE.

- This problem has quite a history:

Reference	Communication complexity	Local gradient complexity	Order	Limitations
DANE [42]	$\mathcal{O}\left(\frac{\delta^2}{\mu^2} \log \frac{1}{\epsilon}\right)$	— <sup>(2)</sup>	1st	quadratic
DiSCo [51]	$\mathcal{O}\left(\sqrt{\frac{\delta}{\mu}} \left(\log \frac{1}{\epsilon} + C^2 \Delta F_0\right) \log \frac{L}{\mu}\right)$	$\mathcal{O}\left(\sqrt{\frac{L}{\mu}} \left(\log \frac{1}{\epsilon} + C^2 \Delta F_0\right) \log \frac{L}{\mu}\right)$	2nd	$C'$ -self-concordant <sup>(1)</sup>
AIDE [40]	$\mathcal{O}\left(\sqrt{\frac{\delta}{\mu}} \log \log \frac{1}{\epsilon}\right)$	$\mathcal{O}\left(\sqrt{\frac{\delta}{\mu}} \sqrt{\frac{L}{\mu}} \log \log \frac{1}{\epsilon}\right)$ <sup>(4)</sup>	1st	quadratic
DANE-LS [50]	$\mathcal{O}\left(\frac{\delta}{\mu} \log \log \frac{1}{\epsilon}\right)$	$\mathcal{O}\left(\sqrt{\frac{L}{\mu}} \frac{\delta^{3/2}}{\sqrt{\mu}} \log \log \frac{1}{\epsilon}\right)$ <sup>(5)</sup>	1st/2nd	quadratic <sup>(6)</sup>
DANE-HB [50]	$\mathcal{O}\left(\sqrt{\frac{\delta}{\mu}} \log \log \frac{1}{\epsilon}\right)$	$\mathcal{O}\left(\sqrt{\frac{L}{\mu}} \frac{\delta}{\mu} \log \log \frac{1}{\epsilon}\right)$ <sup>(5)</sup>	1st/2nd	quadratic <sup>(6)</sup>
SONATA [45]	$\mathcal{O}\left(\frac{\delta}{\mu} \log \log \frac{1}{\epsilon}\right)$	— <sup>(2)</sup>	1st	decentralized
SPAG [21]	$\mathcal{O}\left(\sqrt{\frac{L}{\mu}} \log \log \frac{1}{\epsilon}\right)$ <sup>(1)</sup>	— <sup>(2)</sup>	1st	$M$ -Lipshitz hessian
DiRegINA [12]	$\mathcal{O}\left(\frac{\delta}{\mu} \log \frac{1}{\epsilon} + \sqrt{\frac{M \delta R_0}{\mu}}\right)$	— <sup>(2)</sup>	2nd	$M$ -Lipshitz hessian
ACN [1]	$\mathcal{O}\left(\sqrt{\frac{\delta}{\mu}} \log \log \frac{1}{\epsilon} + \sqrt[3]{\frac{M \delta R_0}{\mu}}\right)$	— <sup>(2)</sup>	2nd	$M$ -Lipshitz hessian
Acc SONATA [46]	$\mathcal{O}\left(\sqrt{\frac{\delta}{\mu}} \log \log \frac{1}{\epsilon} \log \frac{\delta}{\mu}\right)$	— <sup>(2)</sup>	1st	decentralized
This paper	$\mathcal{O}\left(\sqrt{\frac{\delta}{\mu}} \log \log \frac{1}{\epsilon}\right)$	$\mathcal{O}\left(\sqrt{\frac{L}{\mu}} \log \log \frac{1}{\epsilon}\right)$	1st	

In particular, the mirror descent approach with unusual divergence is called DANE.

- The optimal algorithm was proposed in 2022:

Kovalev D. et al. Optimal Gradient Sliding and its Application to Distributed Optimization Under Similarity











Optimal. But may be better?

- **Question:** let's forget about the distribution for 1 slide, and remember is Nesterov's method always optimal? No, if we consider the specificity that the target function can be of the sum species  $f(x) = \frac{1}{n} \sum_{i=1}^n f_i(x)$ .
- **Question:** what then is the optimal method? what are its upper estimates of convergence?

$$f(x) = \frac{1}{n} \sum_{i=1}^n f_i(x).$$

Optimal. But may be better?

- **Question:** let's forget about the distribution for 1 slide, and remember is Nesterov's method always optimal? No, if we consider the specificity that the target function can be of the sum species  $f(x) = \frac{1}{n} \sum_{i=1}^n f_i(x)$ .
- **Question:** what then is the optimal method? what are its upper estimates of convergence?
- The method is called Katyusha, and it has the following upper bound on convergence (oracle complexity by calling  $f_i$ ):

$$\mathcal{O}\left(\left[n + \sqrt{n \frac{L}{\mu}}\right] \log \frac{1}{\varepsilon}\right).$$



Allen-Zhu Z. Katyusha: the first direct acceleration of stochastic gradient methods

**Question:** What is the upper bound on oracular complexity for Nesterov's method?

Optimal. But may be better?

- **Question:** let's forget about the distribution for 1 slide, and remember is Nesterov's method always optimal? No, if we consider the specificity that the target function can be of the sum species  $f(x) = \frac{1}{n} \sum_{i=1}^n f_i(x)$ .
- **Question:** what then is the optimal method? what are its upper estimates of convergence?
- The method is called Katyusha, and it has the following upper bound on convergence (oracle complexity by calling  $f_i$ ):

$$\mathcal{O}\left(\left[n + \sqrt{n \frac{L}{\mu}}\right] \log \frac{1}{\varepsilon}\right).$$



Allen-Zhu Z. Katyusha: the first direct acceleration of stochastic gradient methods

**Question:** What is the upper bound on oracular complexity for Nesterov's method?  $\mathcal{O}\left(n\sqrt{\frac{L}{\mu}} \log \frac{1}{\varepsilon}\right)$

# Variance reduction for similarity

- The idea behind the variance reduction method:

$$\nabla f(x)$$

↓

$$\nabla f_i(x) - \nabla f_i(w) + \nabla f(w),$$

where  $i$  is generated randomly at each iteration from  $[n]$ ,  $w$  – a reference point that is updated infrequently (randomly or deterministically).

- The idea of variance reduction method for data similarity:

$$\nabla f(x) - \nabla f_1(x)$$



$$\nabla f_i(x) - \nabla f_i(w) + \nabla f(w) - f_1(x),$$

where  $i$  is generated randomly at each iteration from  $[M]$ .

## Variance reduction for similarity



Beznosikov A. & Gasnikov A. Compression and data similarity: Combination of two techniques for communication-efficient solving of distributed variational inequalities



Beznosikov A. & Gasnikov A. Similarity, Compression and Local Steps: Three Pillars of Efficient Communications for Distributed Variational Inequalities



Khaled A. & Jin C. Faster federated optimization under second-order similarity

- Previous estimate:

$$\mathcal{O}\left(M\sqrt{1+\frac{\delta}{\mu}\log\frac{1}{\varepsilon}}\right).$$

- We can achieve:

$$\mathcal{O}\left(\left[M + \frac{\delta^2}{\mu^2}\right] \log \frac{1}{\varepsilon}\right) \quad \text{or} \quad \mathcal{O}\left(\left[M + \sqrt{M} \frac{\delta}{\mu}\right] \log \frac{1}{\varepsilon}\right) \quad \text{or}$$

$$\mathcal{O}\left(\left[M + M^{3/4} \sqrt{\frac{\delta}{\mu}}\right] \log \frac{1}{\epsilon}\right)$$