# Progressive & Algorithms & Systems
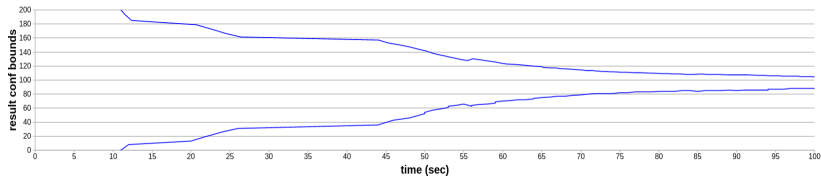
**Florin Rusu**

University of California Merced
Lawrence Berkeley National Laboratory
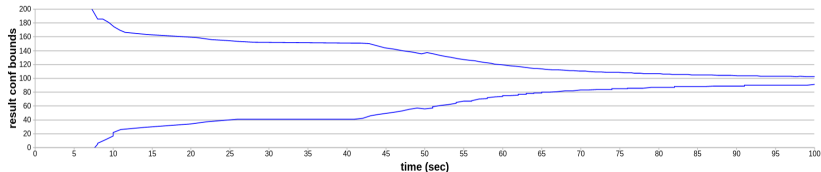
# Progressive Computation for Data Exploration
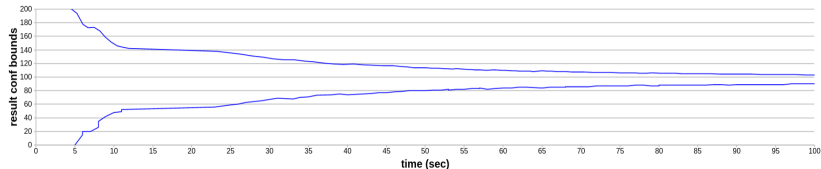


>>> select count(*) from candidates where mag > 10

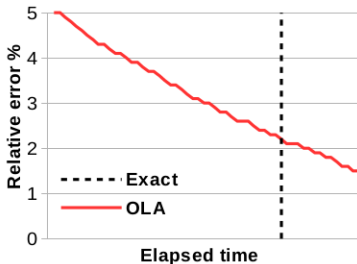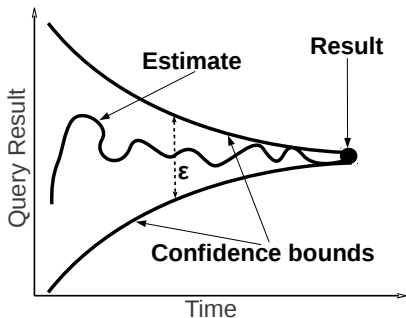>>> select count(*) from candidates where mag > 10 and mag < 100

>>> select count(*) from candidates where mag > 10 and mag < 200
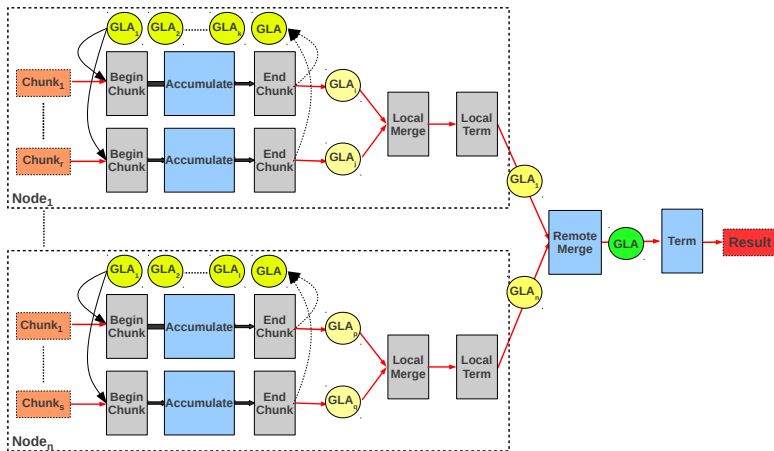
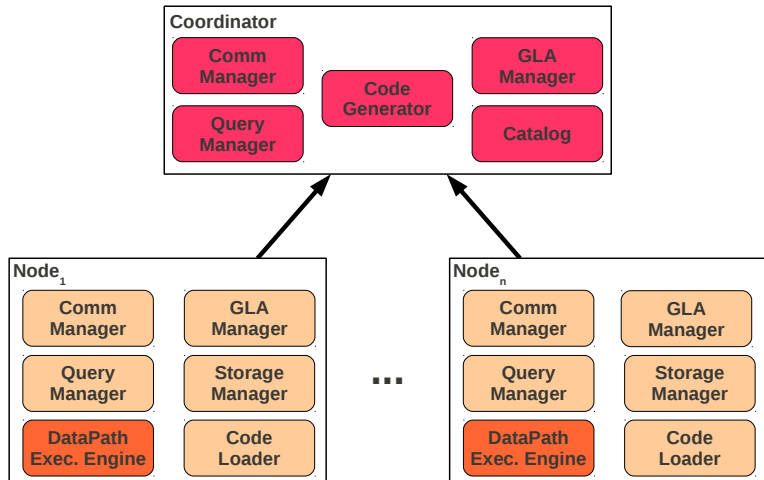# Progressive Computation
## Online Aggregation (OLA) in DB



- How to derive confidence bounds that shrink progressively?
- How to generate a meaningful estimate and confidence bounds as early as possible?
- How to minimize the estimation overhead?

# Outline

# GLADE Architecture

| Method | Usage |
|---|---|
| `Init ()`<br>`Accumulate (Item d)`<br>`Merge (UDA input1,UDA input2, UDA output)` –<br>local and remote<br>`Terminate ()` – local and final | Basic<br>interface |
| `BeginChunk()`<br>`EndChunk ()` | Chunk<br>processing |
| `Serialize ()`<br>`Deserialize ()` | Transfer UDA<br>across processes |
| `EstimatorTerminate ()`<br>`EstimatorMerge (UDA input1,UDA input2, UDA output)` | Progressive<br>computation |
| `Estimate (estimator, lower, upper, confidence)` | OLA estimation |

# Partial Aggregation

# Parallel Sampling

- Centralized random shuffling
  - Permute data randomly at loading
  - Scan produces larger samples

- Stratified sampling
  - Permute data randomly in each partition
  - Direct extension of random shuffling to partitioned data

- Global data randomization at loading
  - Split data randomly at each node
  - Permute all received data randomly
  - Standard hash-based data partitioning

- Centralized
- Distributed tree

# Generic Sampling Estimator

$$\text{AGG} \Longleftarrow \text{SELECT SUM}(f(d))$$
$$\text{FROM } D$$
$$\text{WHERE } P(d)$$

- $S$ is simple random sample without replacement from $D$
- Estimator $X = \frac{|D|}{|S|} \sum_{s \in S, P(s)} f(s)$

$$E[X] = \text{AGG}$$

$$\text{Var}[X] = \frac{|D| - |S|}{(|D| - 1)|S|} \left[ |D| \sum_{d \in D, P(d)} f^2(d) - \left( \sum_{d \in D, P(d)} f(d) \right)^2 \right]$$

$$\text{Est}_{\text{Var}[X]} = \frac{|D|(|D| - |S|)}{|S|^2(|S| - 1)} \left[ |S| \sum_{s \in S, P(s)} f^2(s) - \left( \sum_{s \in S, P(s)} f(s) \right)^2 \right]$$

# Parallel Sampling Estimators

- Data are partitioned across $N$ nodes: $D = D_1 \cup D_2 \cup \cdots \cup D_N$
- Take samples $S_i$, $1 \leq i \leq N$ independently at each node
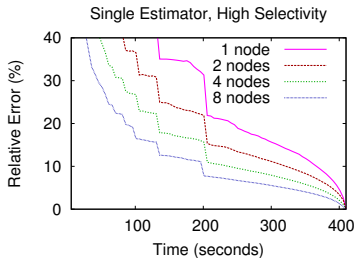
### Single Estimator

- Guarantee
  $S = S_1 \cup S_2 \cup \cdots \cup S_N$ is a sample from $D$

- Synchronized estimator
  - $\frac{S_i}{D_i} = k$ (const),
    $1 \leq i \leq N$

- Asynchronous estimator
  - Global data randomization

### Multiple Estimators

- Stratified sampling
- Build an estimator $X_i$ for each partition $D_i$,
  $1 \leq i \leq N$:
  $X_i = \frac{|D_i|}{|S_i|} \sum_{s \in S_i, P(s)} f(s)$
- $X = \sum_{i=1}^{N} X_i$ is unbiased
- $\mathrm{Var}\left[\sum_{i=1}^{N} X_i\right] = \sum_{i=1}^{N} \mathrm{Var}[X_i]$

# Time to Convergence

```
SELECT n_name, SUM(l_extendprice*(1-l_discount)*(1+l_tax))
FROM lineitem, supplier, nation
WHERE l_shipdate = 1993-02-26 AND l_quantity = 1 AND
 l_discount between [0.02,0.03] AND
 l_suppkey = s_suppkey AND s_nationkey = n_nationkey
GROUP BY n_name
```
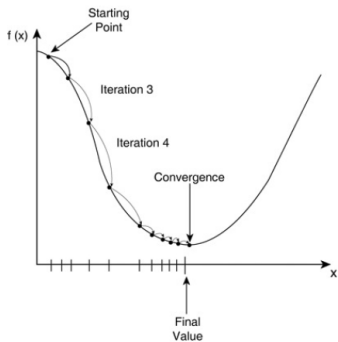


Single Estimator, High Selectivity

- TPC-H scale **8,000 (8TB)**
- Single node: 16 cores @ 2GHz; 16GB RAM; 4 disks @ 110MB/s throughput/disk
- Cluster: 8 X worker + coordinator (9 nodes); Gigabit Ethernet; same rack

| Query | Execution Time (seconds) | |
|---|---|---|
| | No estimation | OLA |
| Aggregate | 222 | 222 |
| $Group_{small}$ | 344 | 345 |
| $Group_{large}$ | 404 | 407 |
| Join | 409 | 411 |

- Chengjie Qin and Florin Rusu. *Sampling Estimators for Parallel Online Aggregation*. BNCOD 2013, pp. 204–217.
- Chengjie Qin and Florin Rusu. *Parallel Online Aggregation in Action*. SSDBM 2013, pp. 383–386. [Demo]
- Chengjie Qin and Florin Rusu. *PF-OLA: A High-Performance Framework for Parallel Online Aggregation*. Distributed and Parallel Databases (DAPD), August 2013.

# Outline

Florin Rusu    Progressive & Algorithms & Systems

# Gradient Descent Optimization



Starting Point

f (x)

Iteration 3

Iteration 4

Convergence

x

Final Value

http://www.yaldex.com/game-development/1592730043_ch18lev1sec4.html

$$min_{\vec{w} \in \mathbb{R}^d} \left\{ \Lambda(\vec{w}) = \sum_{(\vec{x}_i, y_i) \in \text{data}} f\left(\vec{w}, \vec{x}_i; y_i\right) \right\}$$

$$\vec{w}^{(k+1)} = \vec{w}^{(k)} - \alpha^{(k)} \nabla \Lambda\left(\vec{w}^{(k)}\right)$$
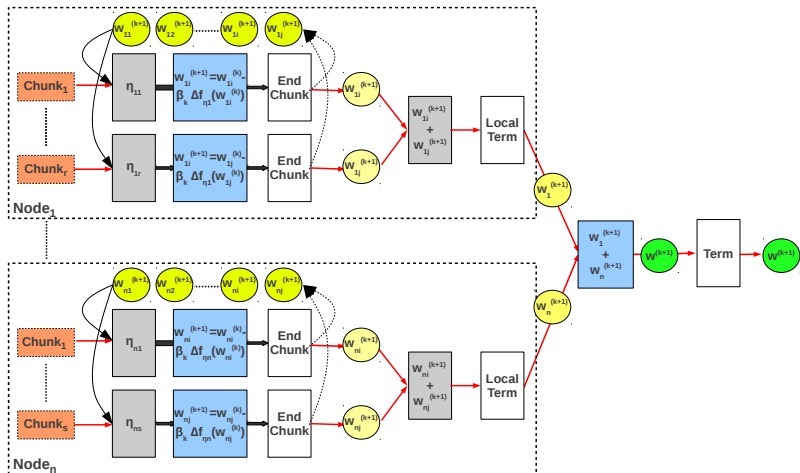
$\Lambda(\vec{w})$ is the loss

$\nabla \Lambda(\vec{w}) = \left[ \dfrac{\partial \Lambda(\vec{w})}{\partial w_1}, \ldots, \dfrac{\partial \Lambda(\vec{w})}{\partial w_d} \right]$ is the gradient

$\alpha^{(k)}$ is step size or learning rate

$\vec{w}^{(0)}$ is the starting point (random)

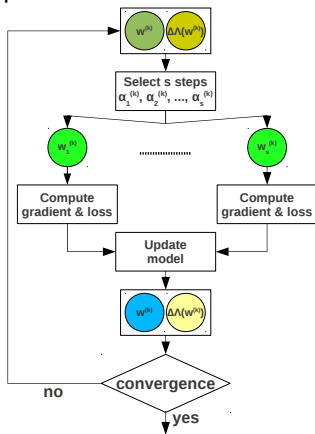- Convergence to minimum guaranteed for convex objective function
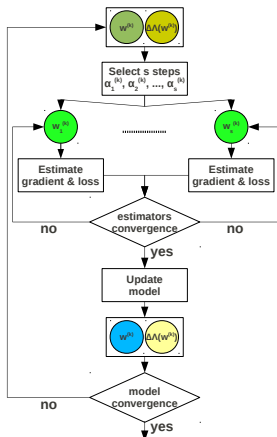
# Gradient Descent as GLADE GLA

# Approximate Gradient Descent

Main idea: apply online aggregation (OLA) sampling to speed-up the execution of a speculative iteration



Speculative Gradient Descent

Approximate Gradient Descent

# Approximate BGD

1: **for each example** $(\vec{x}_{\eta^{(j)}}, y_{\eta^{(j)}})$ **do**
2:   **for each active model** $\vec{w}_i$ **do in parallel**
3:     Estimate gradient $G_i[est, std]$: $\nabla\Lambda(\vec{w}_i)\{(\vec{x}_{\eta^{(j)}}, y_{\eta^{(j)}})\}$
4:     Estimate loss $L_i[est, std]$: $\Lambda(\vec{w}_i, \vec{x}_{\eta^{(j)}}; y_{\eta^{(j)}})$
5:   **end for**
6:   **if** estimators_convergence_check$(j)$ **then**
7:     Prune out models: *Stop Loss*$(\{L_i[est, std]\}_{i \leq s}, \varepsilon_1)$
8:     Let $t$ be the number of remaining models, i.e., *active*
9:     **if** *Stop Gradient*$(\{G_{tl}[est, std]\}_{l \leq d}, \varepsilon_2)$ **then break**
10:     Let $s = t$
11:   **end if**
12: **end for**

```
SELECT
    ..., SUM(f_q), ...
FROM data AS D
SELECT
    ..., Z_q=SUM(f_q), ...
    ..., Z_q^2=SUM(f_q^2), ...
FROM sample AS S
```

- estimator: $Z_q \cdot \frac{|D|}{|S|}$
- variance:
  $\frac{|D|(|D|-|S|)}{|S|^2(|S|-1)}\left(|S|Z_q^2 - (Z_q)^2\right)$

## OLA sampling

- no pre-determined sample size
- avoid correlation between estimators by random data shuffling (at runtime)
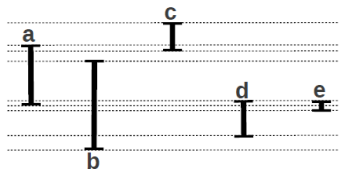- continuous sampling at runtime until estimators converge

# Approximate BGD

1: **for each example** $(\vec{x}_{\eta^{(j)}}, y_{\eta^{(j)}})$ **do**
2:   **for each active model** $\vec{w}_i$ **do in parallel**
3:     Estimate gradient $G_i[est, std]$: $\nabla\Lambda(\vec{w}_i)\{(\vec{x}_{\eta^{(j)}}, y_{\eta^{(j)}})\}$
4:     Estimate loss $L_i[est, std]$: $\Lambda(\vec{w}_i, \vec{x}_{\eta^{(j)}}; y_{\eta^{(j)}})$
5:   **end for**
6:   **if** estimators_convergence_check$(j)$ **then**
7:     Prune out models: *Stop Loss*$(\{L_i[est, std]\}_{i \le s}, \varepsilon_1)$
8:     Let $t$ be the number of remaining models, i.e., *active*
9:     **if** *Stop Gradient*$(\{G_{tl}[est, std]\}_{l \le d}, \varepsilon_2)$ **then break**
10:     Let $s = t$
11:   **end if**
12: **end for**

- Identify step size with minimum loss based on estimate & confidence bounds
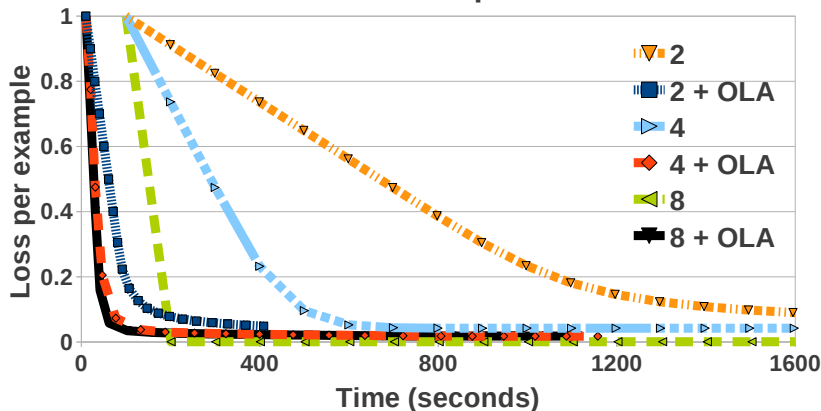- Interacting estimators
- Exact pruning
- Approximate pruning

# Approximate BGD

1: **for each example** $(\vec{x}_{\eta^{(j)}}, y_{\eta^{(j)}})$ **do**

2:      **for each active model** $\vec{w}_i$ **do in parallel**

3:          Estimate gradient $G_i[est, std]$: $\nabla \Lambda(\vec{w}_i)\{(\vec{x}_{\eta^{(j)}}, y_{\eta^{(j)}})\}$

4:          Estimate loss $L_i[est, std]$: $\Lambda(\vec{w}_i, \vec{x}_{\eta^{(j)}}; y_{\eta^{(j)}})$

5:      **end for**

6:      **if** estimators_convergence_check($j$) **then**

7:          Prune out models: *Stop Loss*$(\{L_i[est, std]\}_{i \le s}, \varepsilon_1)$

8:          Let $t$ be the number of remaining models, i.e., *active*

9:          **if** *Stop Gradient*$(\{G_{tl}[est, std]\}_{l \le d}, \varepsilon_2)$ **then break**

10:          Let $s = t$

11:      **end if**

12: **end for**

- One estimator for each dimension in model $\vec{w}$
- Grouped estimators

**When to stop?**

- all estimators converge (million dimensions)
- percentage of estimators converge, e.g., $90\%$
- unified convergence threshold, e.g., $\varepsilon' = d \cdot \varepsilon$

BGD OLA # step sizes

- 50M examples, 200 dimensions, 136GB
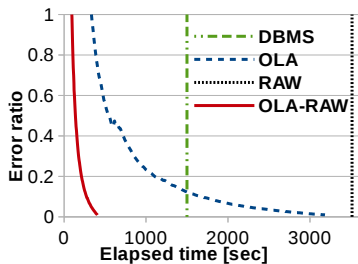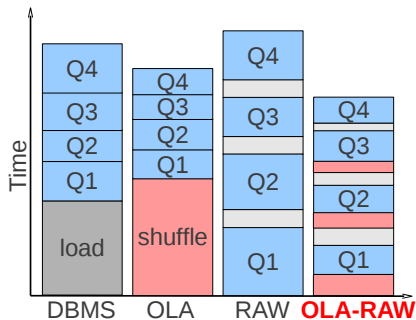
# Sample Percentage
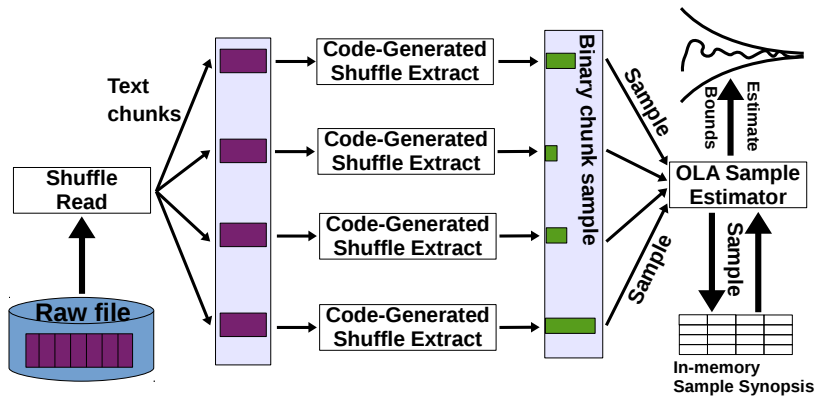


- 50M examples, 200 dimensions, 136GB

- Chengjie Qin and Florin Rusu. *Speculative Approximations for Terascale Analytics*. CoRR abs/1501.00255, December 2014.
- Chengjie Qin and Florin Rusu. *Speculative Approximations for Terascale Distributed Gradient Descent Optimization*. DanaC@SIGMOD 2015.
- Florin Rusu, Chengjie Qin, and Martin Torres. *Scalable Analytics Model Calibration with Online Aggregation*. IEEE Data Engineering Bulletin, Vol. 38, No. 3, pp. 30–44, September 2015.

# Outline

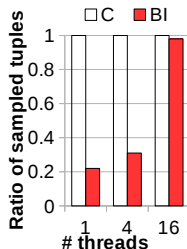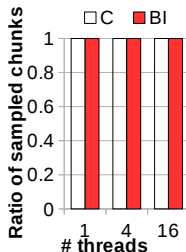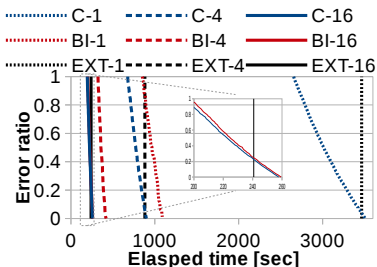Florin Rusu    Progressive & Algorithms & Systems

- Inspection paradox: result order $\neq$ random chunk order

# Experimental Results



Selectivity = 100%

Selectivity = 50%

- Yu Cheng, Weijie Zhao, and Florin Rusu. *OLA-RAW: Scalable Exploration over Raw Data*. CoRR abs/1702.00358, February 2017.
- Yu Cheng, Weijie Zhao, and Florin Rusu. *Bi-Level Online Aggregation on Raw Data*. SSDBM 2017.

# Thank you!

# Questions?