

# AWARE: Workload-aware, Redundancy-exploiting Linear Algebra

Sebastian Baunsgaard & Matthias Boehm

## Lossy

- ❖ Main stream ML Systems
- ❖ Specialized data types
- ❖ Associated with trial and error

## Lossy

- ❖ Main stream ML Systems
- ❖ Specialized data types
- ❖ Associated with trial and error

## Lossless CLA

## Lossy

- ❖ Main stream ML Systems
- ❖ Specialized data types
- ❖ Associated with trial and error

## Lossless CLA

- ❖ Lightweight database compression
- ❖ Preaggregation Cocoding
- ❖ Correct results

## Lossy

- ❖ Main stream ML Systems
- ❖ Specialized data types
- ❖ Associated with trial and error

## Lossless CLA

- ❖ Lightweight database compression
- ❖ Preaggregation Cocoding
- ❖ Correct results

## Optimization Goal

- ❖ Normal: **Compression Ratio**
- ❖ We: **Workload Time**

## SystemDS



<https://github.com/apache/systemds>

## Reproducibility



<https://github.com/damslab/reproducibility>

## User Script:

```
X = read("data/X")  
y = read("data/y")  
  
X = scale(X, TRUE, TRUE)  
w = l2svm(X, y, TRUE,  
          1e-9, 1e-3, 100)  
  
write(w, "data/wXy")
```

## User Script:

```
X = read("data/X")
y = read("data/y")

X = scale(X, TRUE, TRUE)
w = l2svm(X, y, TRUE,
          1e-9, 1e-3, 100)

write(w, "data/wXy")
```

## Built-in Functions:

```
if(shift)
  X = X - colMeans(X)
if(scale)
  X = X / colSds(X)
```

```
if(intercept)
  X = cbind(X, ones)
while(conto & i < maxi) {
  Xd = X %*% s
  while(conti) {
    out = 1 - y * (Xw + sz * Xd)
    sz = sz - g/h; # ...
  }
  g_new = t(X) %*% (out * y)
}
```



## User Script:

```
X = read("data/X")  
y = read("data/y")
```

```
X = scale(X, TRUE, TRUE)  
w = l2svm(X, y, TRUE,  
          1e-9, 1e-3, 100)
```

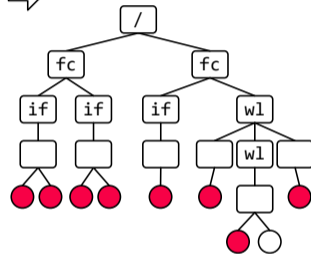
```
write(w, "data/wXy")
```

## Built-in Functions:

```
if(shift)  
  X = X - colMeans(X)  
if(scale)  
  X = X / colSds(X)
```

```
if(intercept)  
  X = cbind(X, ones)  
while(conto & i < maxi) {  
  Xd = X %*% s  
  while(conti) {  
    out = 1 - y * (Xw + sz * Xd)  
    sz = sz - g/h; # ...  
  }  
  g_new = t(X) %*% (out * y)  
}
```

## Workload Tree



## User Script:

```
X = read("data/X")  
y = read("data/y")
```

```
X = scale(X, TRUE, TRUE)  
w = l2svm(X, y, TRUE,  
          1e-9, 1e-3, 100)
```

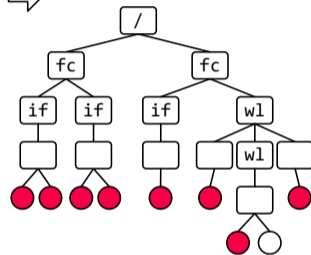
```
write(w, "data/wXy")
```

## Built-in Functions:

```
if(shift)  
  X = X - colMeans(X)  
if(scale)  
  X = X / colSds(X)
```

```
if(intercept)  
  X = cbind(X, ones)  
while(conto & i < maxi) {  
  Xd = X %*% s  
  while(conti) {  
    out = 1 - y * (Xw + sz * Xd)  
    sz = sz - g/h; # ...  
  }  
  g_new = t(X) %*% (out * y)  
}
```

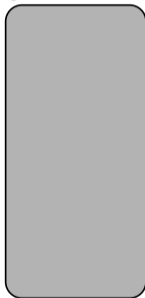
## Workload Tree

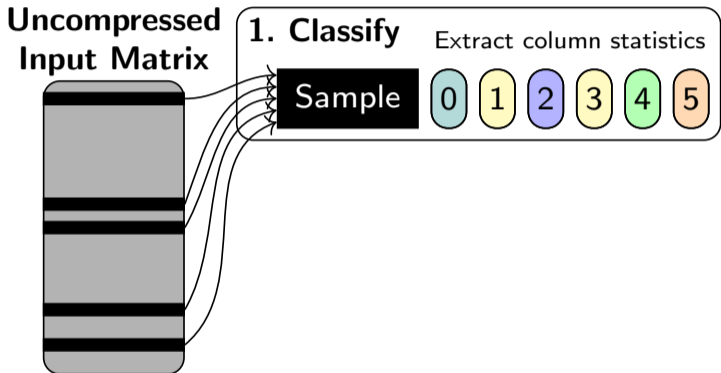


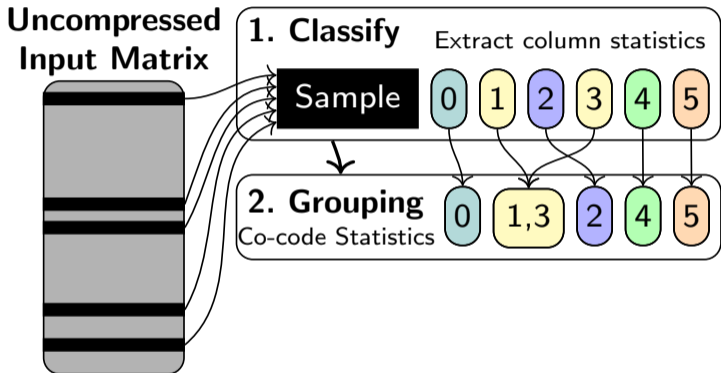
## Cost Summary

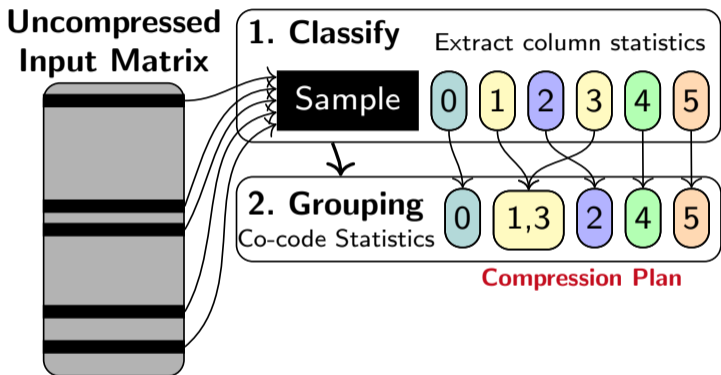
0	100	10	10	105	0
---	-----	----	----	-----	---

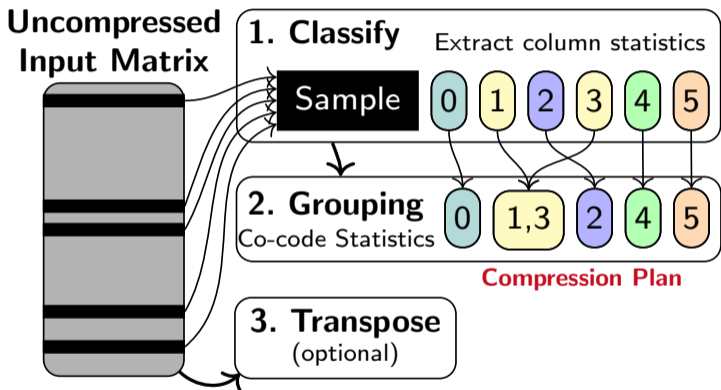
Uncompressed  
Input Matrix

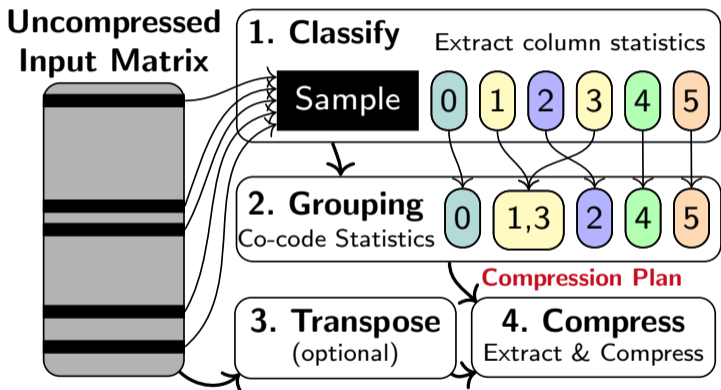




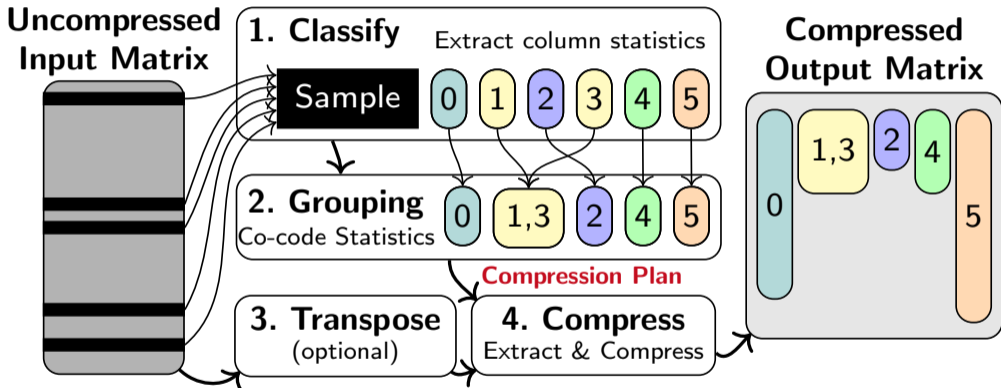












# Compression Example



**Uncompressed**

1	2	0	6	0	4
7	3	0	4	0	7
1	2	0	6	0	4
7	2	0	5	0	7
1	3	9	4	0	4
7	0	8.5	0	0	7
1	3	8.5	4	0	7
7	3	0	4	0	7
1	2	0	6	0	4
7	3	9	4	0	7

# Compression Example

## Uncompressed

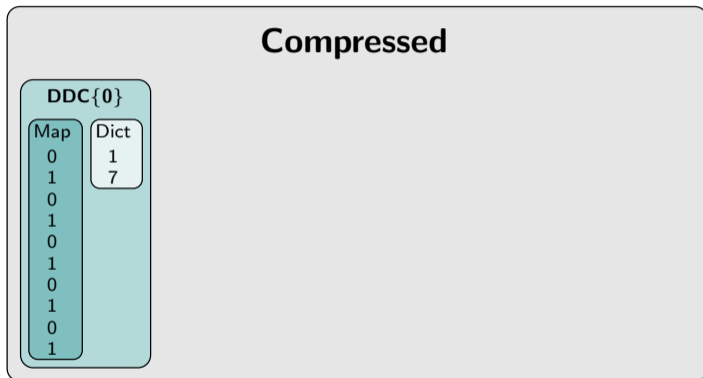
1	2	0	6	0	4
7	3	0	4	0	7
1	2	0	6	0	4
7	2	0	5	0	7
1	3	9	4	0	4
7	0	8.5	0	0	7
1	3	8.5	4	0	7
7	3	0	4	0	7
1	2	0	6	0	4
7	3	9	4	0	7

## Compressed

# Compression Example

### Uncompressed

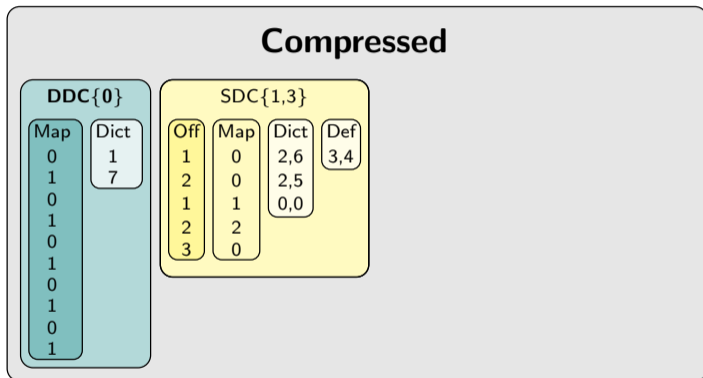
1	2	0	6	0	4
7	3	0	4	0	7
1	2	0	6	0	4
7	2	0	5	0	7
1	3	9	4	0	4
7	0	8.5	0	0	7
1	3	8.5	4	0	7
7	3	0	4	0	7
1	2	0	6	0	4
7	3	9	4	0	7



# Compression Example

### Uncompressed

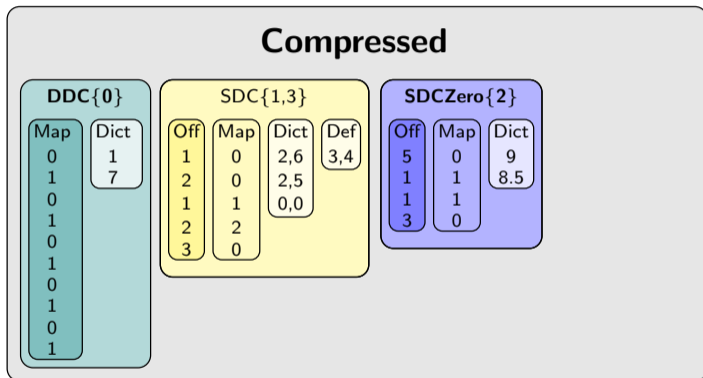
1	2	0	6	0	4
7	3	0	4	0	7
1	2	0	6	0	4
7	2	0	5	0	7
1	3	9	4	0	4
7	0	8.5	0	0	7
1	3	8.5	4	0	7
7	3	0	4	0	7
1	2	0	6	0	4
7	3	9	4	0	7



# Compression Example

### Uncompressed

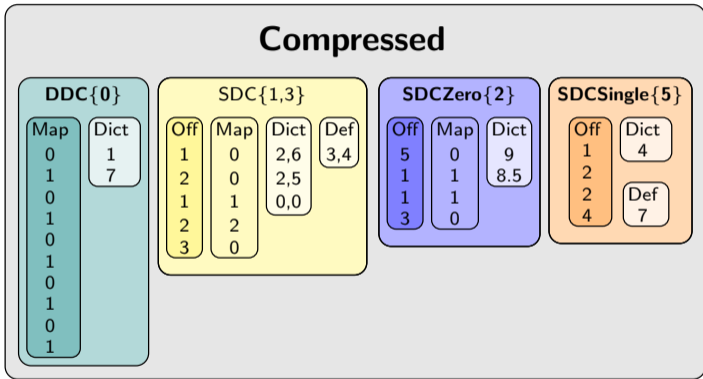
1	2	0	6	0	4
7	3	0	4	0	7
1	2	0	6	0	4
7	2	0	5	0	7
1	3	9	4	0	4
7	0	8.5	0	0	7
1	3	8.5	4	0	7
7	3	0	4	0	7
1	2	0	6	0	4
7	3	9	4	0	7



# Compression Example

**Uncompressed**

1	2	0	6	0	4
7	3	0	4	0	7
1	2	0	6	0	4
7	2	0	5	0	7
1	3	9	4	0	4
7	0	8.5	0	0	7
1	3	8.5	4	0	7
7	3	0	4	0	7
1	2	0	6	0	4
7	3	9	4	0	7



	CLA	Aware
Co-Coding	$\mathcal{O}(m^2)$	
Column Group Encodings	4 (5)	
Materialization	Eager	
Optimization Objective	Data	
Matrix Multiplication	MV, VM	



	CLA	Aware
Co-Coding	$\mathcal{O}(m^2)$	$\mathcal{O}(m)$
Column Group Encodings	4 (5)	
Materialization	Eager	
Optimization Objective	Data	
Matrix Multiplication	MV, VM	

	CLA	Aware
Co-Coding	$\mathcal{O}(m^2)$	$\mathcal{O}(m)$
Column Group Encodings	4 (5)	7 (327)
Materialization	Eager	
Optimization Objective	Data	
Matrix Multiplication	MV, VM	

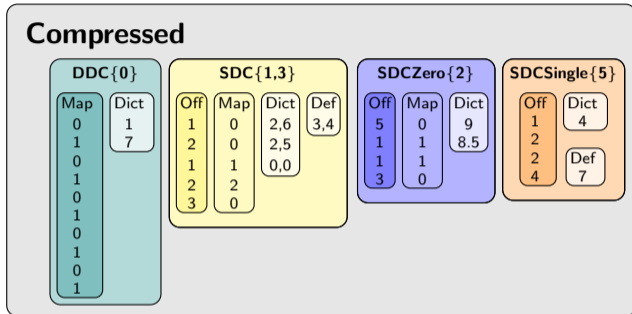
	CLA	Aware
Co-Coding	$\mathcal{O}(m^2)$	$\mathcal{O}(m)$
Column Group Encodings	4 (5)	7 (327)
Materialization	Eager	Deferred
Optimization Objective	Data	
Matrix Multiplication	MV, VM	

	CLA	Aware
Co-Coding	$\mathcal{O}(m^2)$	$\mathcal{O}(m)$
Column Group Encodings	4 (5)	7 (327)
Materialization	Eager	Deferred
Optimization Objective	Data	Data & Ops
Matrix Multiplication	MV, VM	

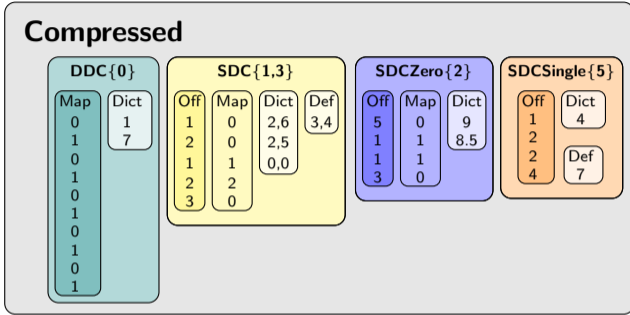
	CLA	Aware
Co-Coding	$\mathcal{O}(m^2)$	$\mathcal{O}(m)$
Column Group Encodings	4 (5)	7 (327)
Materialization	Eager	Deferred
Optimization Objective	Data	Data & Ops
Matrix Multiplication	MV, VM	MV, VM, MM

Type	Description	CLA	Aware
CON	Constant or Empty Columns		✓
DDC	Dense Dictionary Coding	✓	✓
OLE	Offset-list Encoding	✓	(✓)
FOR	Frame of Reference		✓
RLE	Run-length Encoding	✓	(✓)
SDC	Sparse Dictionary Coding		✓
UC	Uncompressed (dense/sparse)	✓	✓

# Operation Example



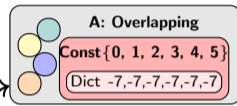
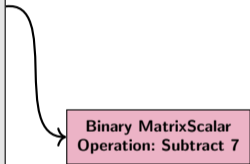
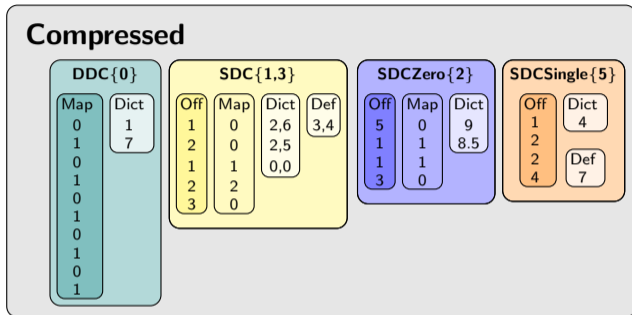
# Operation Example



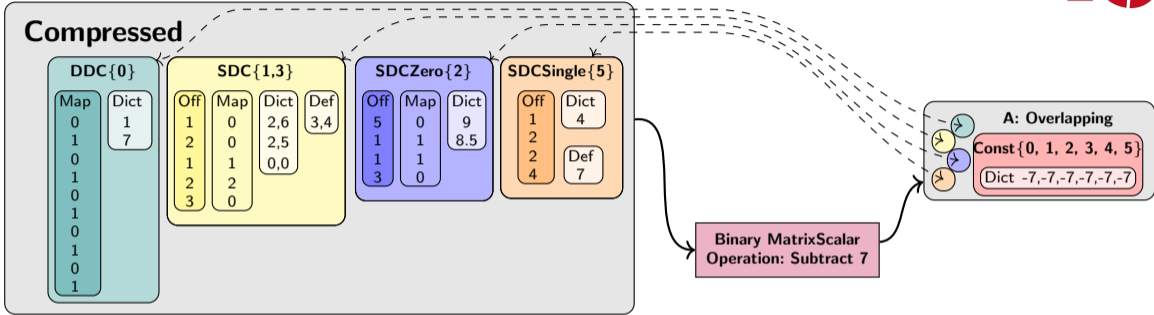
Binary MatrixScalar  
Operation: Subtract 7



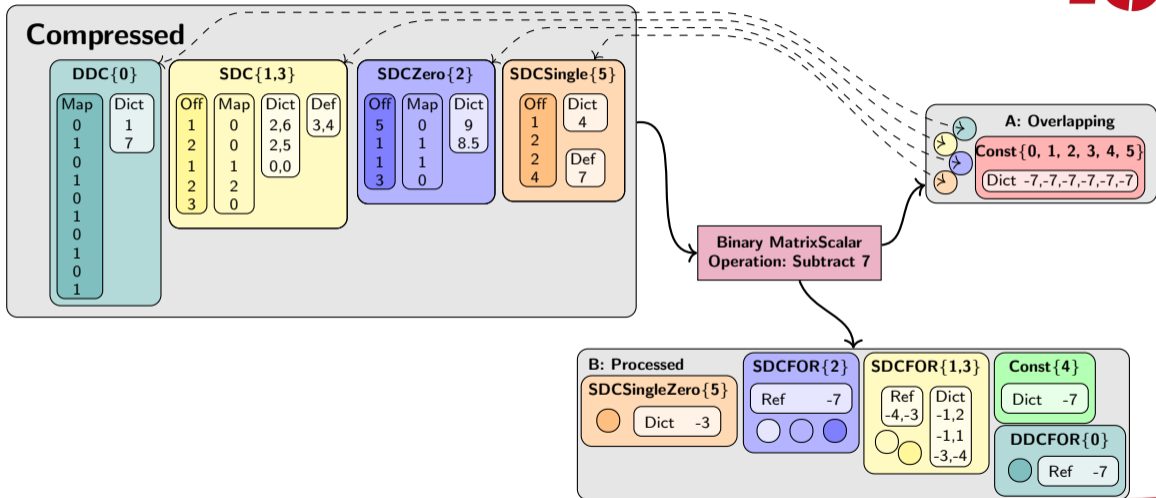
# Operation Example



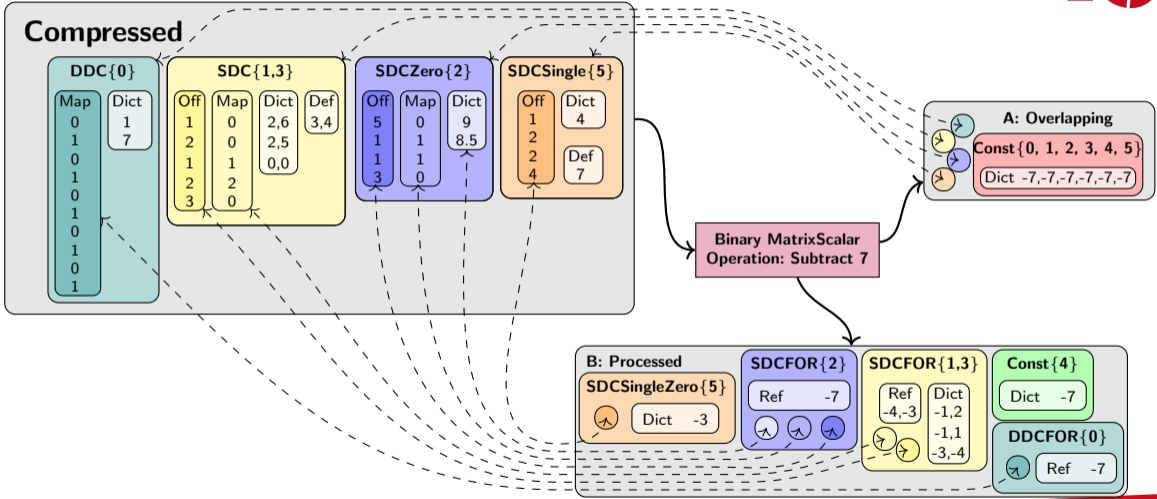
# Operation Example



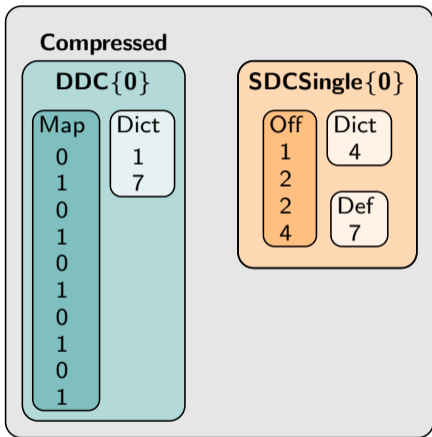
# Operation Example



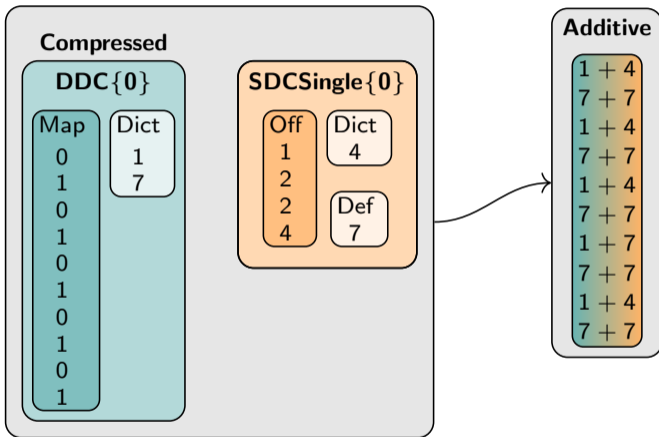
# Operation Example



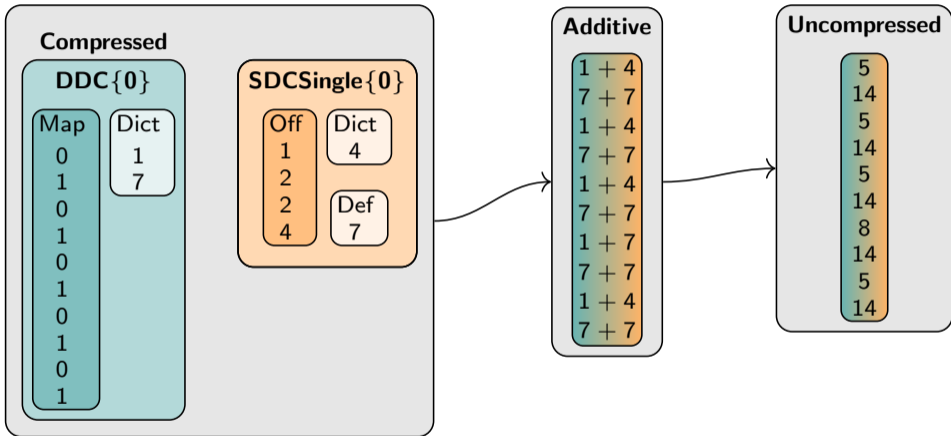
# Overlapping Example



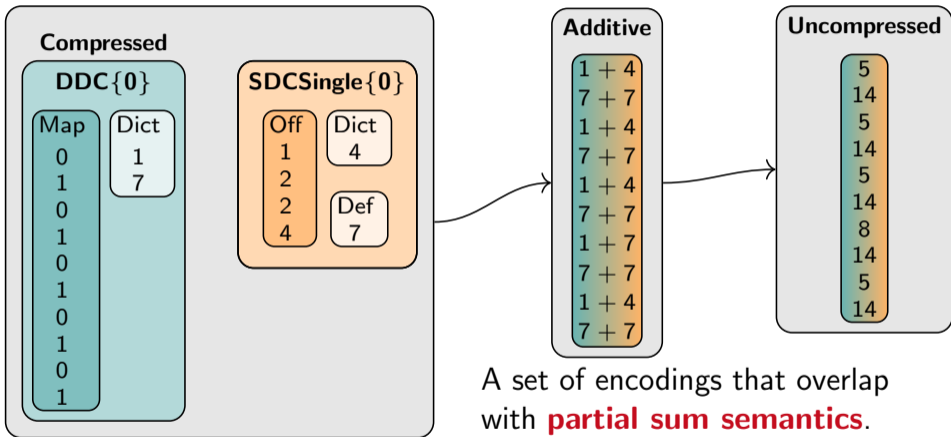
# Overlapping Example



# Overlapping Example

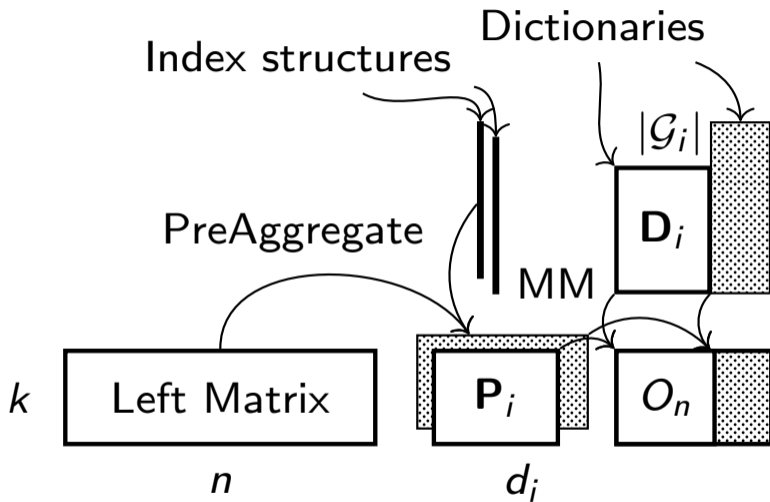


# Overlapping Example

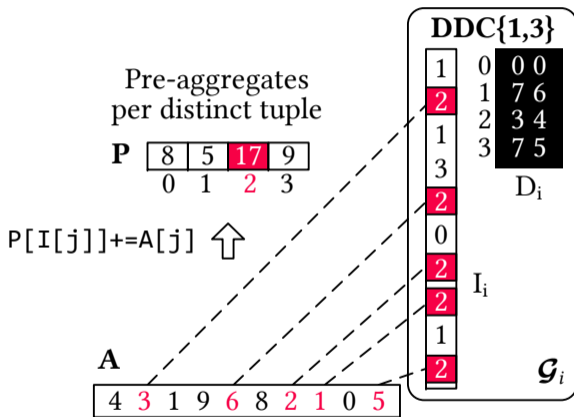




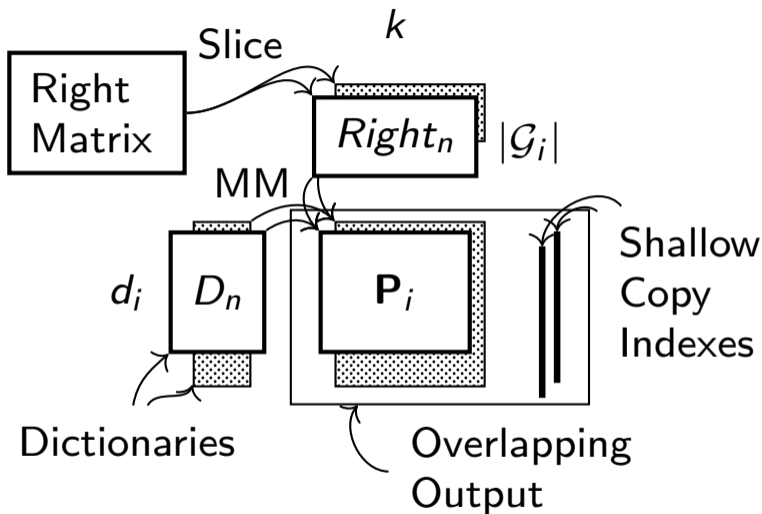
# Left Matrix Multiplication



# Preaggregate



# Right Matrix Multiplication



Datasets ( $n$  Rows,  $m$  Columns,  $sp$  Sparsity).

Dataset	$n$ (nrow( $\mathbf{X}$ ))	$m$ (ncol( $\mathbf{X}$ ))	$sp$	Size
Airline78	14,462,943	29	0.54	3.4 GB
Amazon	8,026,324	<b>2,330,066</b>	1.2e-6	1.22 GB
Covtype	581,012	54	0.22	127 MB
Mnist1m	1,000,000	784	0.25	2.46 GB
<b>US Census</b>	2,458,285	68 <b>(378)</b>	0.43 (0.18)	1.34 GB
US Census 128x	314,660,480	68 (378)	0.43 (0.18)	289.5 GB

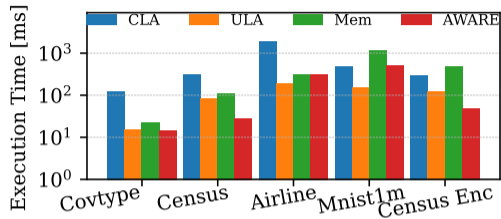
**Million column** dataset and **high cocoding potential** scale up dataset

Local Compression Times [Seconds] and Ratios.

Dataset	CLA		Aware-Mem		Aware	
	time	ratio	time	ratio	time	ratio
Airline78	9.34 sec	10.22	1.74 sec	8.61	2.08 sec	7.94
Amazon	<i>37.6 hours</i>	<i>Crash</i>	8.54 sec	1.73	3.77 sec	<i>Abort</i>
Covtype	1.10 sec	13.79	0.84 sec	14.24	1.23 sec	13.99
Mnist1m	7.25 sec	7.14	4.57 sec	6.09	<i>17.50 sec</i>	4.41
US Census	5.15 sec	35.38	1.16 sec	29.60	1.15 sec	27.35
US Census Enc	<i>27.48 sec</i>	41.03	<i>5.78 sec</i>	38.46	6.54 sec	29.46

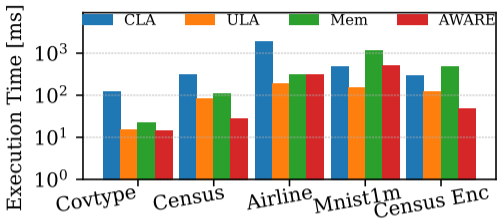
**Intelligent compression abort** based on workload and  
**Faster compression**

# Matrix Multiplication

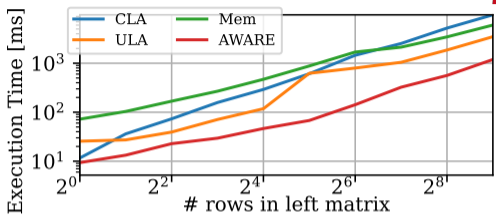


(a) 16 row LMM

# Matrix Multiplication

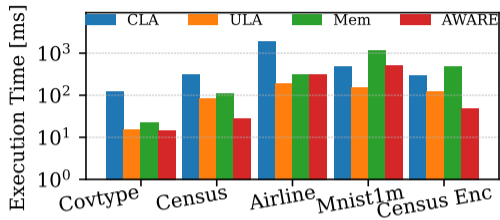


(e) 16 row LMM

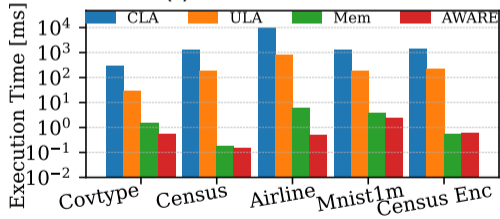


(f) LMM Census Enc Scaling

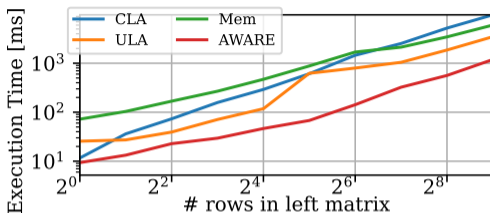
# Matrix Multiplication



(i) 16 row LMM



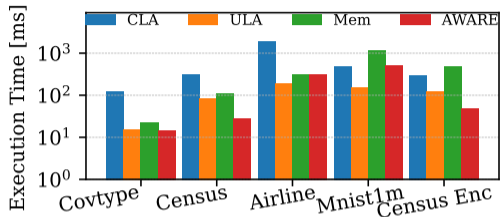
(k) 16 col RMM



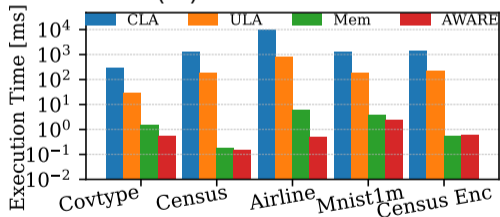
(j) LMM Census Enc Scaling



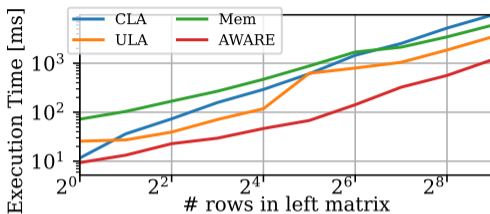
# Matrix Multiplication



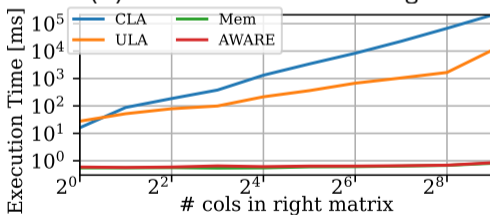
(m) 16 row LMM



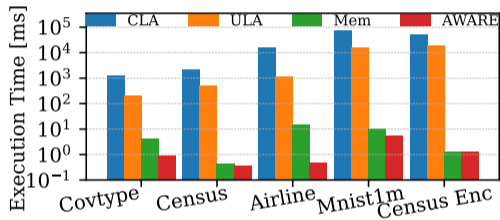
(o) 16 col RMM



(n) LMM Census Enc Scaling

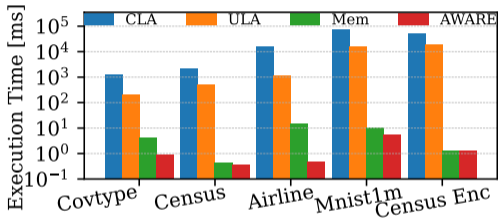


(p) RMM Census Enc Scaling

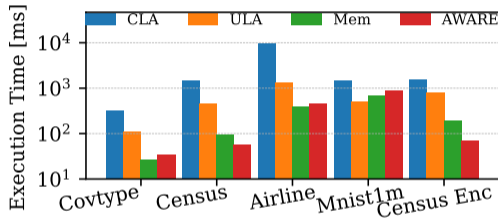


(q) Scale and Shift Sparse

**Operation sequences** maintain **compressed formats** and



(s) Scale and Shift Sparse



(t) Euclidean MinDist 16 Points Dense

**Operation sequences** maintain **compressed formats** and **Decompressing sequences** are competitive or **faster** than uncompressed

# Local End-to-End Algorithms

Workload-awareness on Local End-to-End Algorithms (Data: US Census Enc)

	ULA	Aware-Mem		Aware	
	Time	Comp	Time	Comp	Time
<b>K-Means</b>	51.6 sec	4.2 sec	46.2 sec	6.2 sec	27.1 sec
<b>PCA</b>	12.7 sec	4.0 sec	10.4 sec	6.0 sec	9.0 sec
<b>MLogReg</b>	32.0 sec	4.5 sec	32.5 sec	7.2 sec	26.0 sec
<b>ImCG</b>	19.8 sec	5.0 sec	20.7 sec	6.4 sec	18.6 sec
<b>ImDS</b>	15.6 sec	5.7 sec	15.5 sec	6.1 sec	14.3 sec
<b>L2SVM</b>	38.9 sec	6.5 sec	45.2 sec	6.2 sec	36.5 sec

**AWARE** is able to **amortize** the cost of **online compression** in short jobs.

# End-to-End Algorithms Hybrid Execution



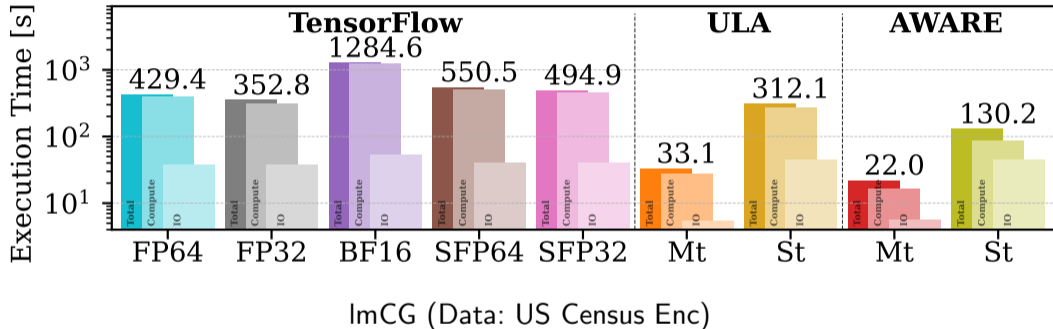
[Seconds] (Data: US Census Enc,  $D$  .. Incl. Distributed Ops).

	K-Means		PCA		MLogReg	
	ULA	Aware	ULA	Aware	ULA	Aware
1x	51.6	(6) 27.1	12.7	(6) 9.4	32.0	(7) 26.0
8x	471.0	(26) 117.8	330.3	(26) 42.6	393.3	(29) 88.2
16x	$D$ 484.3	(48) 183.9	$D$ 76.3	(47) 67.5	$D$ 570.3	(58) 144.2
32x	$D$ 1,491.6	$D$ 1,496.3	$D$ 70.3	$D$ 61.2	$D$ 671.5	$D$ 629.9
128x	$D$ 17,819.0	$D$ 6,298.0	$D$ 137.0	$D$ 140.3	$D$ 3,502.9	$D$ 1,710.6
*128x	$D$ 33,039.0	$D$ 11,616.0	$D$ 269.0	$D$ 259.0	$D$ 50,998.0	$D$ 8,599.6

**AWARE** keeps operations local, since they **fit in memory**.



# TensorFlow Comparison



Our **baseline is competitive**, and **AWARE is the fastest**.

GridSearch MLogReg (Data: US Census Enc).

ULA	Aware-Mem	Aware
274.3 sec	238.1 sec	<b>92.6 sec</b>

Aware reduce **memory bandwidth** and **reduce compute** requirements.

- ❖ Optimize for **execution cost** not size of data
- ❖ Use Compressed Linear Algebra for **redundancy exploitation**, it is the natural next step from **sparsity exploitation**
- ❖ Workload awareness exploit **workload and data** characteristics together



- ❖ Optimize for **execution cost** not size of data
- ❖ Use Compressed Linear Algebra for **redundancy exploitation**, it is the natural next step from **sparsity exploitation**
- ❖ Workload awareness exploit **workload and data** characteristics together

## Future Work

- ❖ Compressed **feature transformation** of heterogeneous data
- ❖ Orthogonal **lossy compression**