

# Fine-grained Provenance for Linear Algebra Operators

Zhepeng Yan

Val Tannen

Zachary Ives

University of Pennsylvania

# Motivation

- Provenance is well-understood for **relational data / queries**.
  - E.g., view maintenance, delete propagation, computing trust, prob. db
- But increasingly analysts are performing more complex tasks:
  - Machine learning, data mining, image analysis, graph analytics
- **Array data** and **matrix algebra** are commonly used!

# Motivation

- Provenance is well-understood for **relational data / queries**.
  - E.g., view maintenance, delete propagation, computing trust, prob. db
- But increasingly analysts are performing more complex tasks:
  - Machine learning, data mining, image analysis, graph analytics
- **Array data** and **matrix algebra** are commonly used!
- Question: How do we **track provenance** in this setting?

# Inspiration: Provenance Semirings

- An algebra framework [Green et al. PODS'07] for
  - **annotating** tuples in a relation
  - **propagating** annotations through relational queries (SPJU and aggregation)
- Enables efficient **delete propagation, view maintenance, etc**

# Semiring example: input data

## P

<b>PID</b>	<b>PValue</b>
101	11
102	12
2003	13
2004	14
2005	15
2006	16

## Q

<b>QID</b>	<b>QValue</b>
2003	0
2004	0
2005	0
2006	0

## R

<b>RID</b>	<b>RValue</b>
101	13
102	14
5005	15
5006	16
5007	17
5008	18
5009	19

# Semiring example: query

**P**

PID	PValue
101	11
102	12
2003	13
2004	14
2005	15
2006	16

**Q**

QID	QValue
2003	0
2004	0
2005	0
2006	0

**R**

RID	RValue
101	13
102	14
5005	15
5006	16
5007	17
5008	18
5009	19

**S(Value) :- P(x, Value), R(x, \_)**

**S(Value) :- P(x, Value), Q(x, \_)**

**S(Value) :- R(\_, Value)**

# Semiring example: query

## P

PID	PValue
101	11
102	12
2003	13
2004	14
2005	15
2006	16

## Q

QID	QValue
2003	0
2004	0
2005	0
2006	0

## R

RID	RValue
101	13
102	14
5005	15
5006	16
5007	17
5008	18
5009	19

**S(Value) :- P(x, Value), R(x, \_)**

S(Value) :- P(x, Value), Q(x, \_)

S(Value) :- R(\_, Value)

# Semiring example: query

## P

PID	PValue
101	11
102	12
2003	13
2004	14
2005	15
2006	16

## Q

QID	QValue
2003	0
2004	0
2005	0
2006	0

## R

RID	RValue
101	13
102	14
5005	15
5006	16
5007	17
5008	18
5009	19

$S(\text{Value}) :- P(x, \text{Value}), R(x, \_)$

**$S(\text{Value}) :- P(x, \text{Value}), Q(x, \_)$**

$S(\text{Value}) :- R(\_, \text{Value})$



# Semiring example: query

## P

PID	PValue
101	11
102	12
2003	13
2004	14
2005	15
2006	16

## Q

QID	QValue
2003	0
2004	0
2005	0
2006	0

## R

RID	RValue
101	13
102	14
5005	15
5006	16
5007	17
5008	18
5009	19

$S(\text{Value}) :- P(x, \text{Value}), R(x, \_)$

$S(\text{Value}) :- P(x, \text{Value}), Q(x, \_)$

**$S(\text{Value}) :- R(\_, \text{Value})$**

# Semiring example: output tuples

**S**

Value
11
12
13
14
15
16
17
18
19

**S(Value) :- P(x, Value), R(x, \_)****S(Value) :- P(x, Value), Q(x, \_)****S(Value) :- R(\_, Value)**

# Semiring example: annotated output

**S**

**S(Value) :- P(x, Value), R(x, \_)**

**S(Value) :- P(x, Value), Q(x, \_)**

**S(Value) :- R(\_, Value)**

Value	Annotation
11	<i>pr</i>
12	<i>pr</i>
13	<i>pq + r</i>
14	<i>pq + r</i>
15	<i>pq + r</i>
16	<i>pq + r</i>
17	<i>r</i>
18	<i>r</i>
19	<i>r</i>

# Semiring example: delete propagation

## S

Value	Annotation
11	$pr$
12	$pr$
13	$pq + r$
14	$pq + r$
15	$pq + r$
16	$pq + r$
17	$r$
18	$r$
19	$r$

What if we remove tuples from **P**?

Set  $p = 0$ !

# Semiring example: delete propagation

## S

Value	Annotation
11	<i>0</i>
12	<i>0</i>
13	<i>r</i>
14	<i>r</i>
15	<i>r</i>
16	<i>r</i>
17	<i>r</i>
18	<i>r</i>
19	<i>r</i>

What if we remove tuples from **P**?

Set  $p = 0$ !

# Semimodule example: aggregation

**S****Query: SUM(Value)**

Value	Annotation
11	$pr$
12	$pr$
13	$pq + r$
14	$pq + r$
15	$pq + r$
16	$pq + r$
17	$r$
18	$r$
19	$r$

# Semimodule example: annotation

**S**

**Query: SUM(Value)**

Value	Annotation
11	$pr$
12	$pr$
13	$pq + r$
14	$pq + r$
15	$pq + r$
16	$pq + r$
17	$r$
18	$r$
19	$r$

**Annotated aggregation**

$$pr * (11 + 12) + (pq + r) * (13 + 14 + 15 + 16) + r * (17 + 18 + 19)$$

# Semimodule example: annotation

**S**

**Query: SUM(Value)**

Value	Annotation
11	$pr$
12	$pr$
13	$pq + r$
14	$pq + r$
15	$pq + r$
16	$pq + r$
17	$r$
18	$r$
19	$r$

**Annotated aggregation**

$$pr * (11 + 12) + (pq + r) * (13 + 14 + 15 + 16) + r * (17 + 18 + 19)$$

**First term: annotation**



# Semimodule example: annotation

**S**

**Query: SUM(Value)**

Value	Annotation
11	$pr$
12	$pr$
13	$pq + r$
14	$pq + r$
15	$pq + r$
16	$pq + r$
17	$r$
18	$r$
19	$r$

**Annotated aggregation**

$$pr * (11 + 12) + (pq + r) * (13 + 14 + 15 + 16) + r * (17 + 18 + 19)$$

First term: annotation

**Second term: value**

# Semimodule example: annotation

**S**

**Query: SUM(Value)**

Value	Annotation
11	$pr$
12	$pr$
13	$pq + r$
14	$pq + r$
15	$pq + r$
16	$pq + r$
17	$r$
18	$r$
19	$r$

**Annotated aggregation**

$$pr * (11 + 12) + (pq + r) * (13 + 14 + 15 + 16) + r * (17 + 18 + 19)$$

or  $pr * 23 + (pq + r) * 58 + r * 54$

# Semimodule example: annotation

**S**

**Query: SUM(Value)**

Value	Annotation
11	$pr$
12	$pr$
13	$pq + r$
14	$pq + r$
15	$pq + r$
16	$pq + r$
17	$r$
18	$r$
19	$r$

**Annotated aggregation**

$$pr * (11 + 12) + (pq + r) * (13 + 14 + 15 + 16) + r * (17 + 18 + 19)$$

$$\text{or } pr * 23 + (pq + r) * 58 + r * 54$$

$$\text{or } pq * 58 + r * (p * 23 + 112)$$

# Semimodule example: annotation

## S

Query: SUM(Value)

Value	Annotation
11	$pr$
12	$pr$
13	$pq + r$
14	$pq + r$
15	$pq + r$
16	$pq + r$
17	$r$
18	$r$
19	$r$

**Annotated aggregation**

$$pr * (11 + 12) + (pq + r) * (13 + 14 + 15 + 16) + r * (17 + 18 + 19)$$

or  $pr * 23 + (pq + r) * 58 + r * 54$

or  $pq * 58 + r * (p * 23 + 112)$

**Delete propagation:**

set  $r = 0$  and obtain  $pq * 58$

# Tracking matrix provenance

- We want to get the same benefits!
  - Algebraically manipulate annotated matrices
  - Hypothetical deletion

# Tracking matrix data

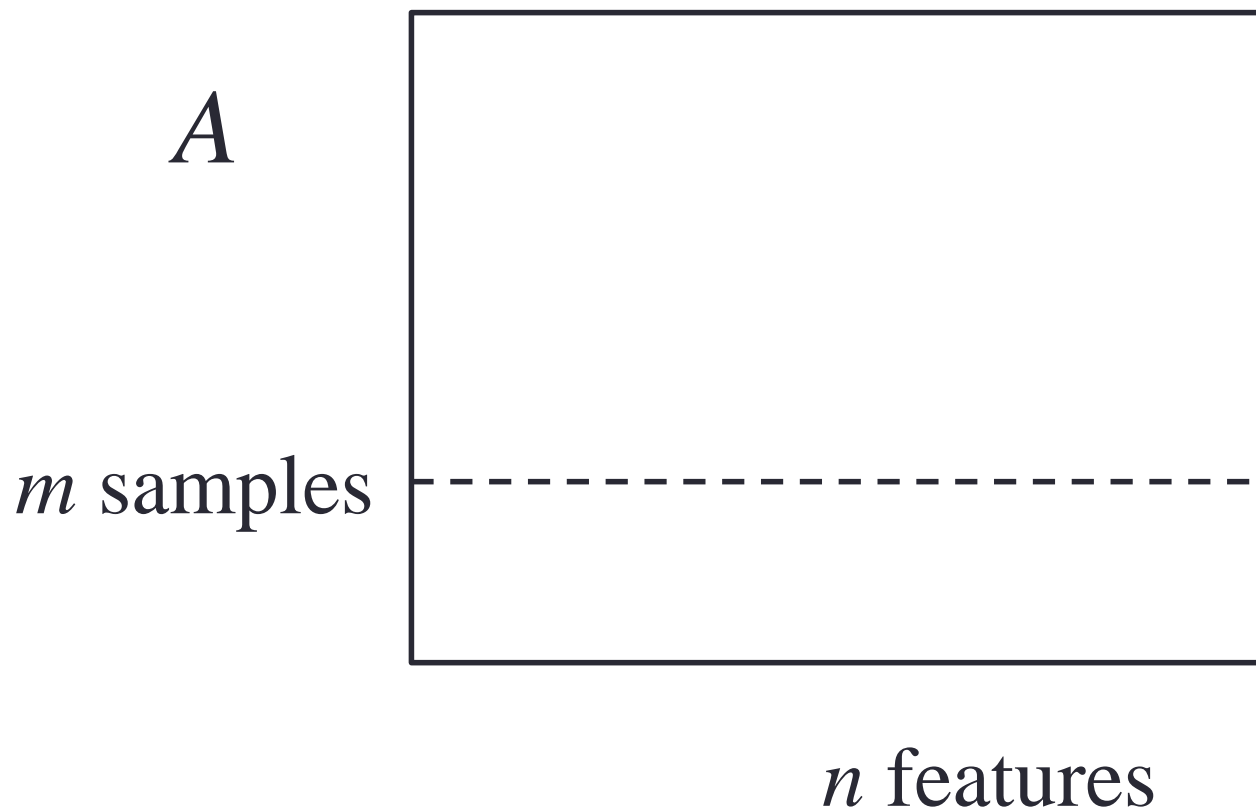
$A$

$m$  samples

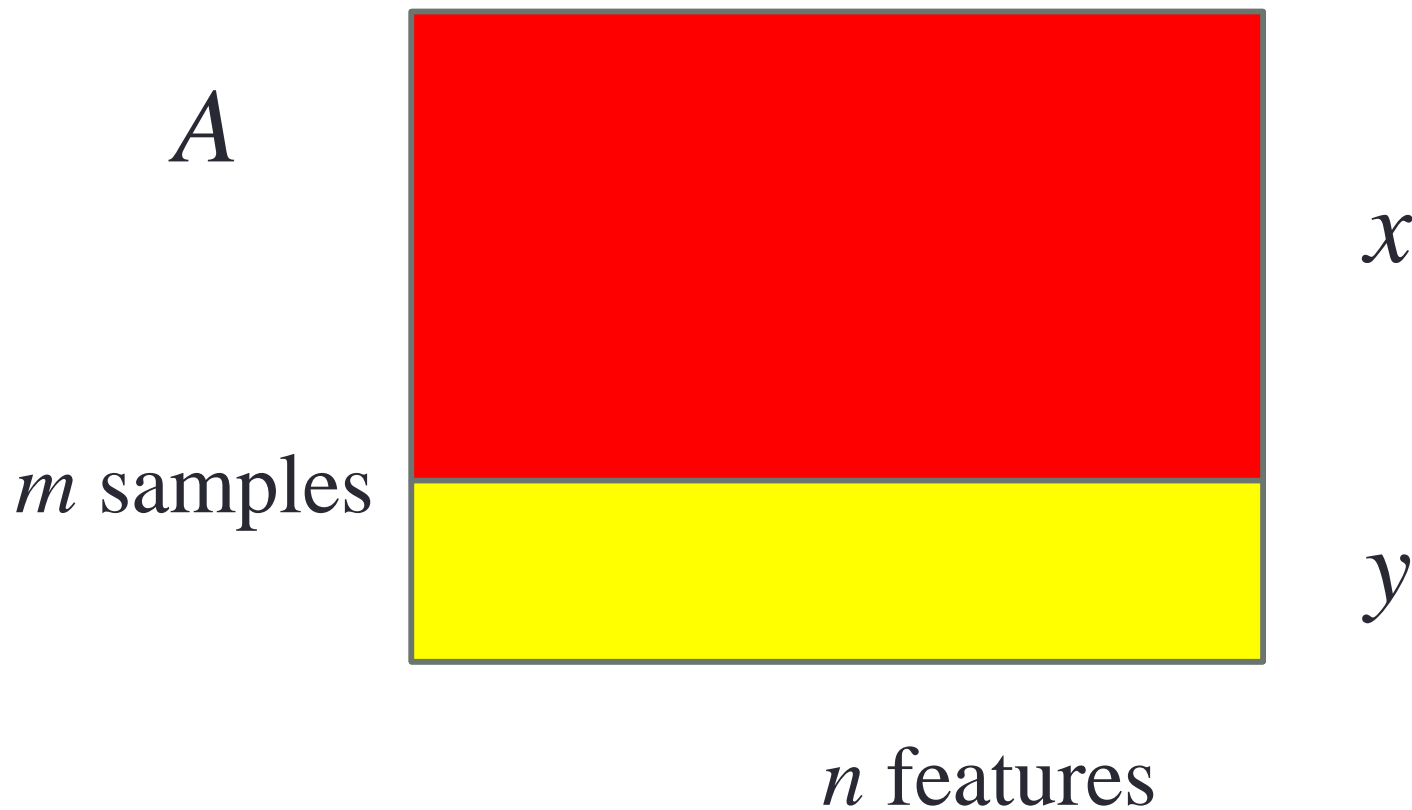


$n$  features

# Tracking matrix data: partitioning

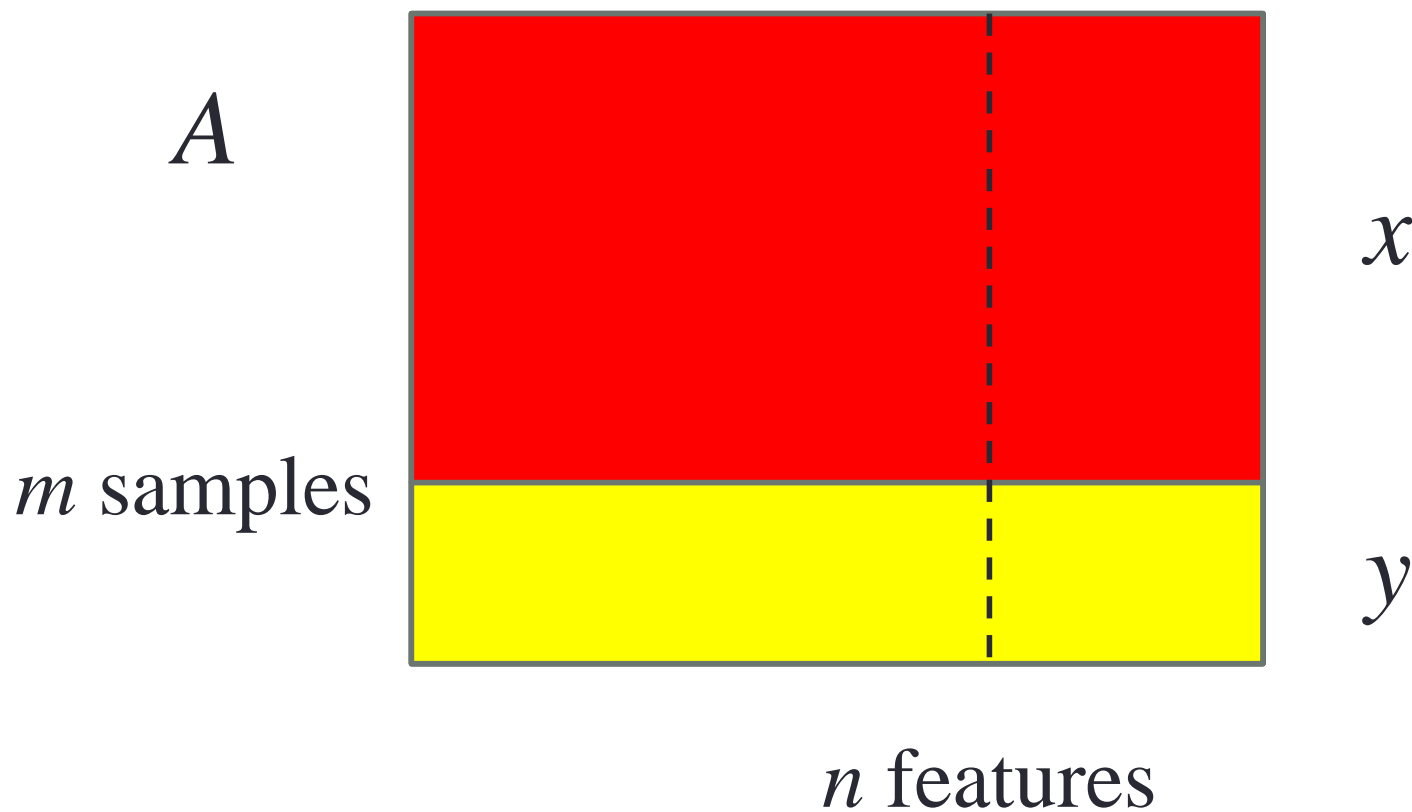


# Tracking matrix data: annotating

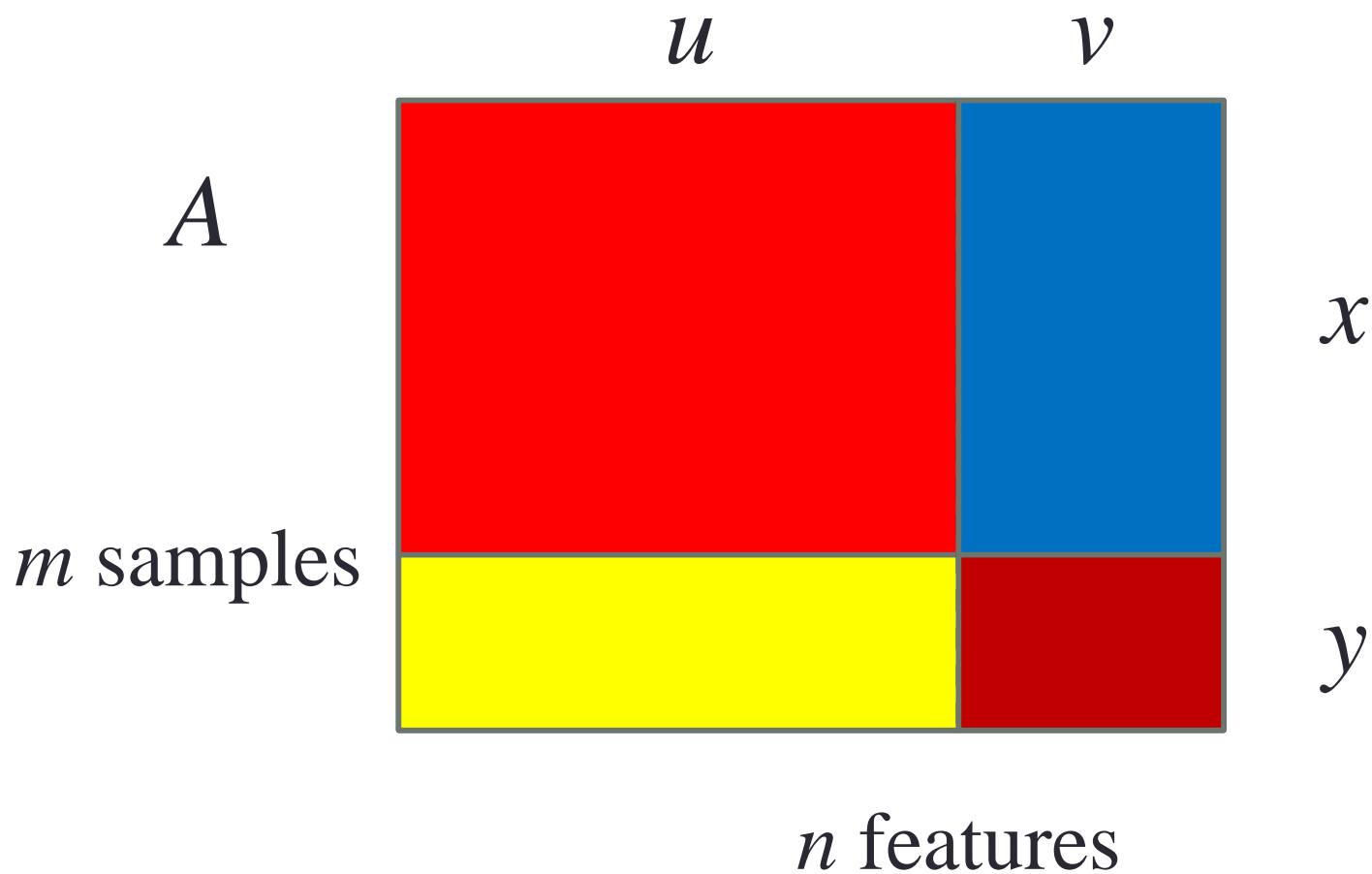




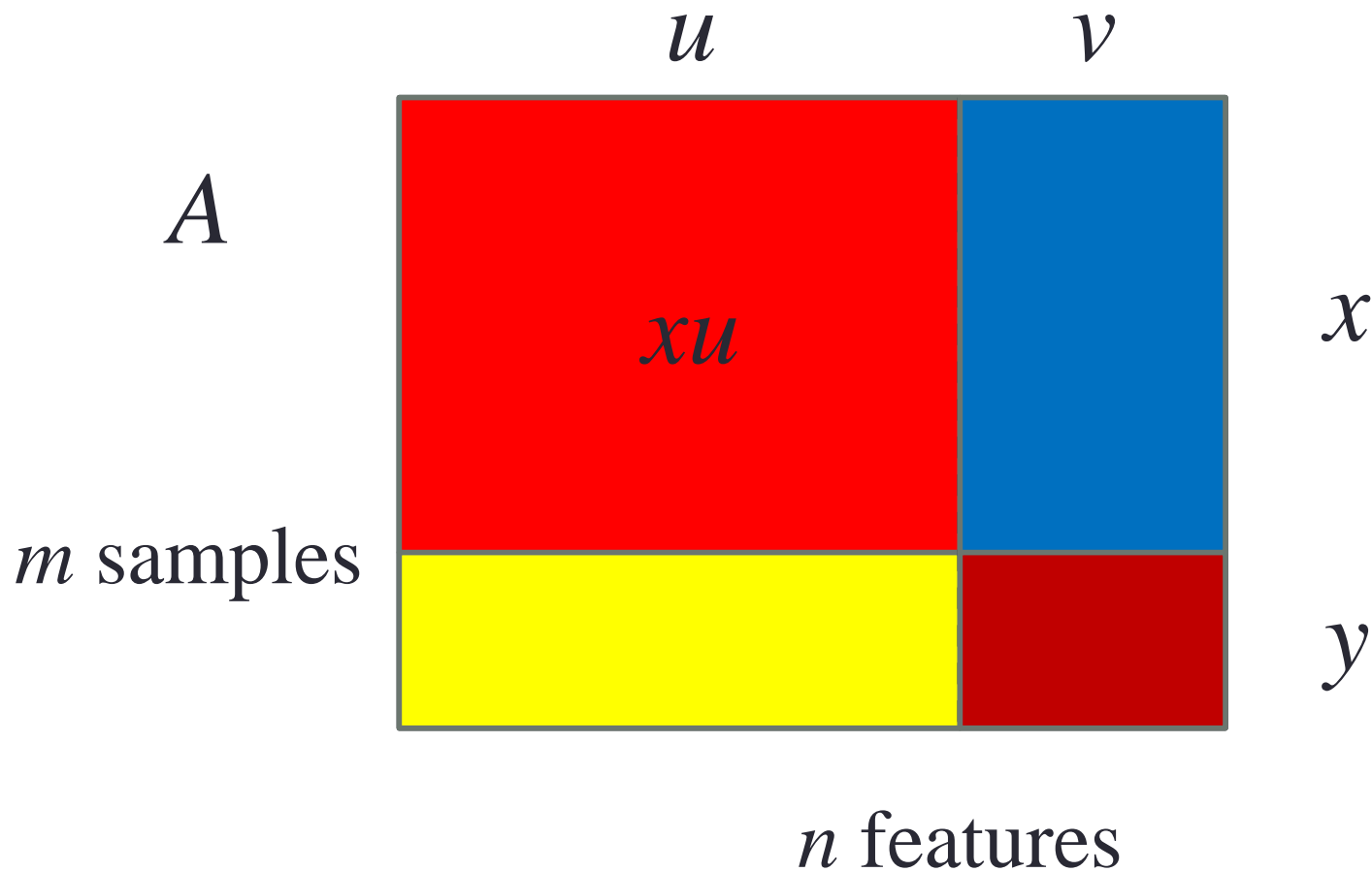
# Tracking matrix data: annotating



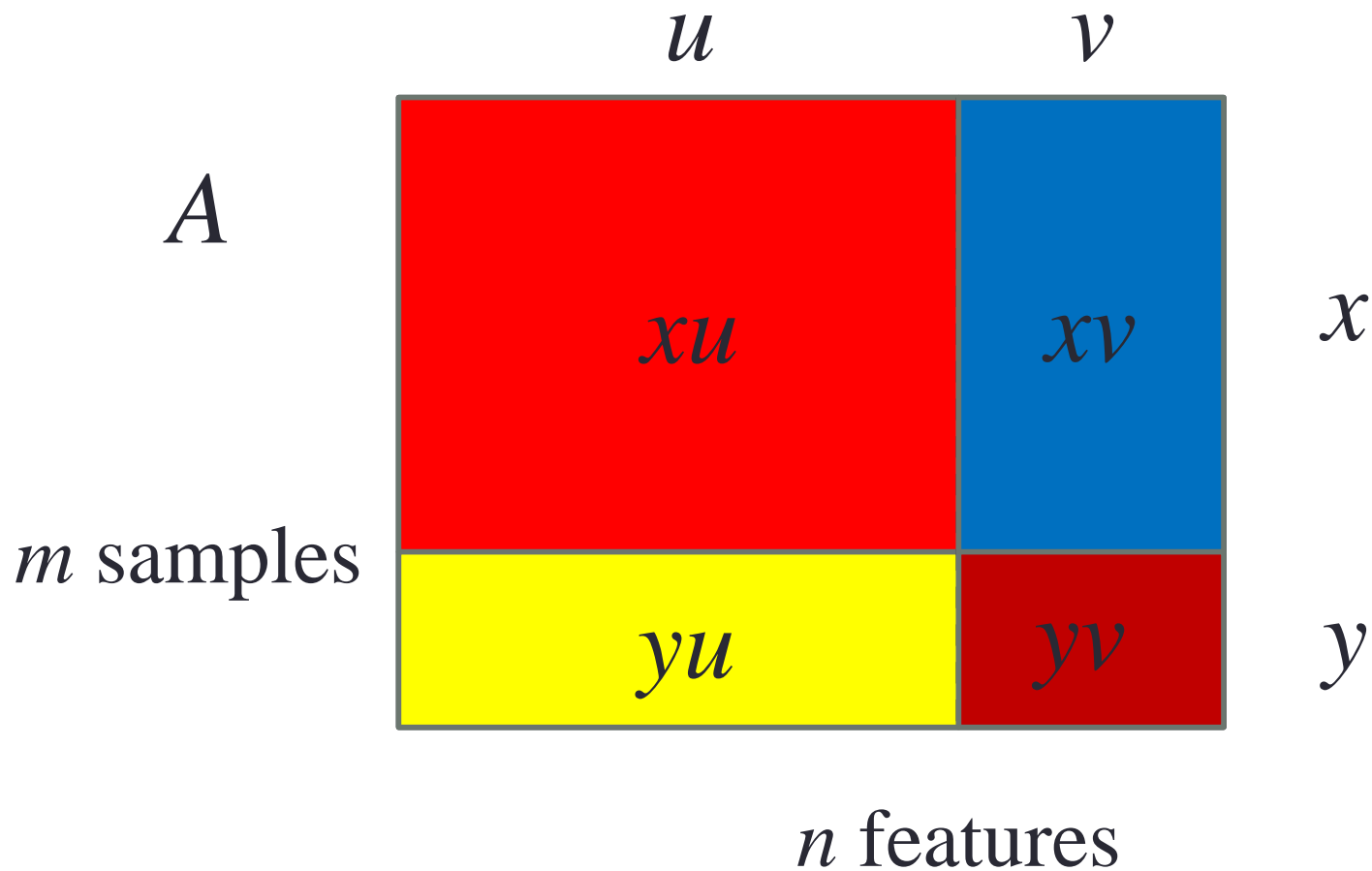
# Tracking matrix data: annotating



# Tracking matrix data: annotating



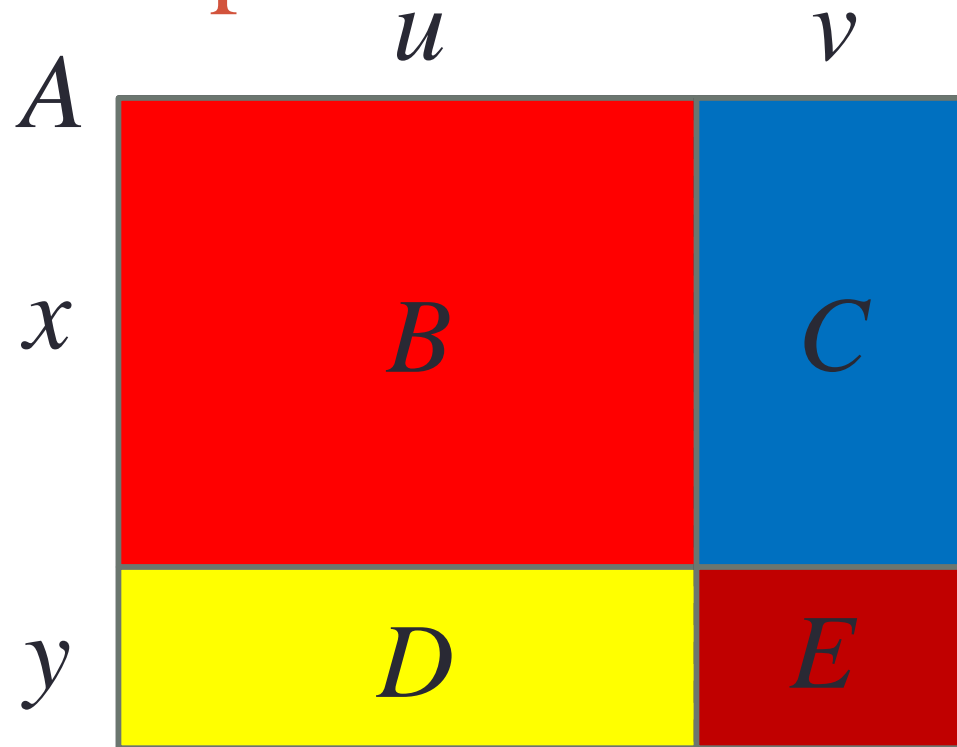
# Tracking matrix data: annotating



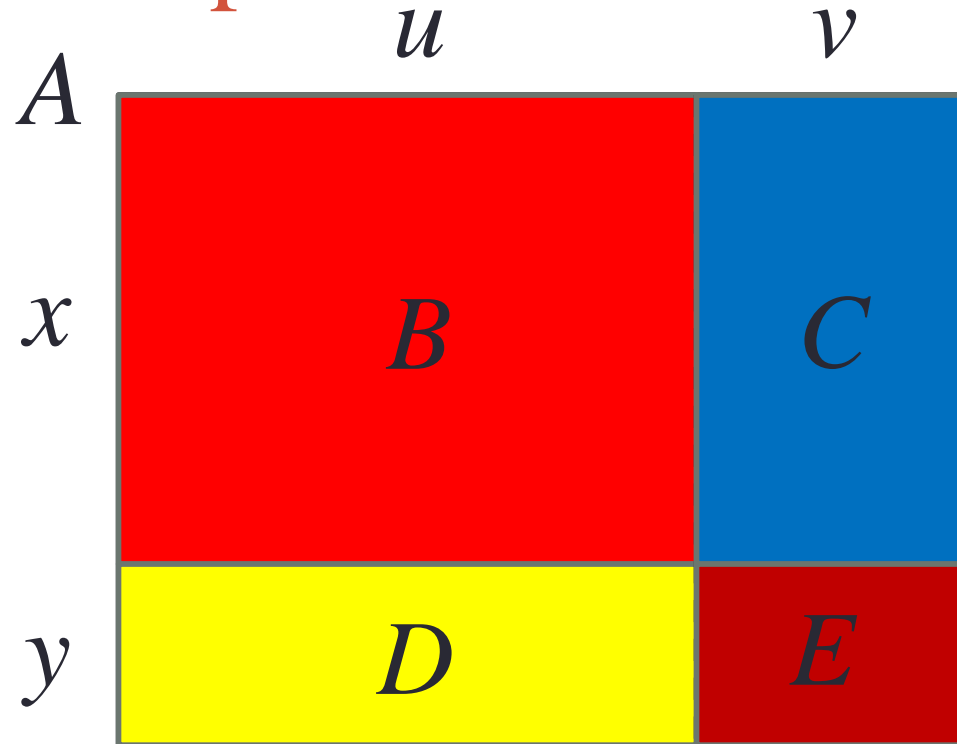
# Challenges

- Specify and relate different parts of a given matrix
  - Matrix decomposition through **selector matrices**
- Specify connection between derived and source matrices
  - Embed matrix algebra and provenance into a **semialgebra**

# Decomposition



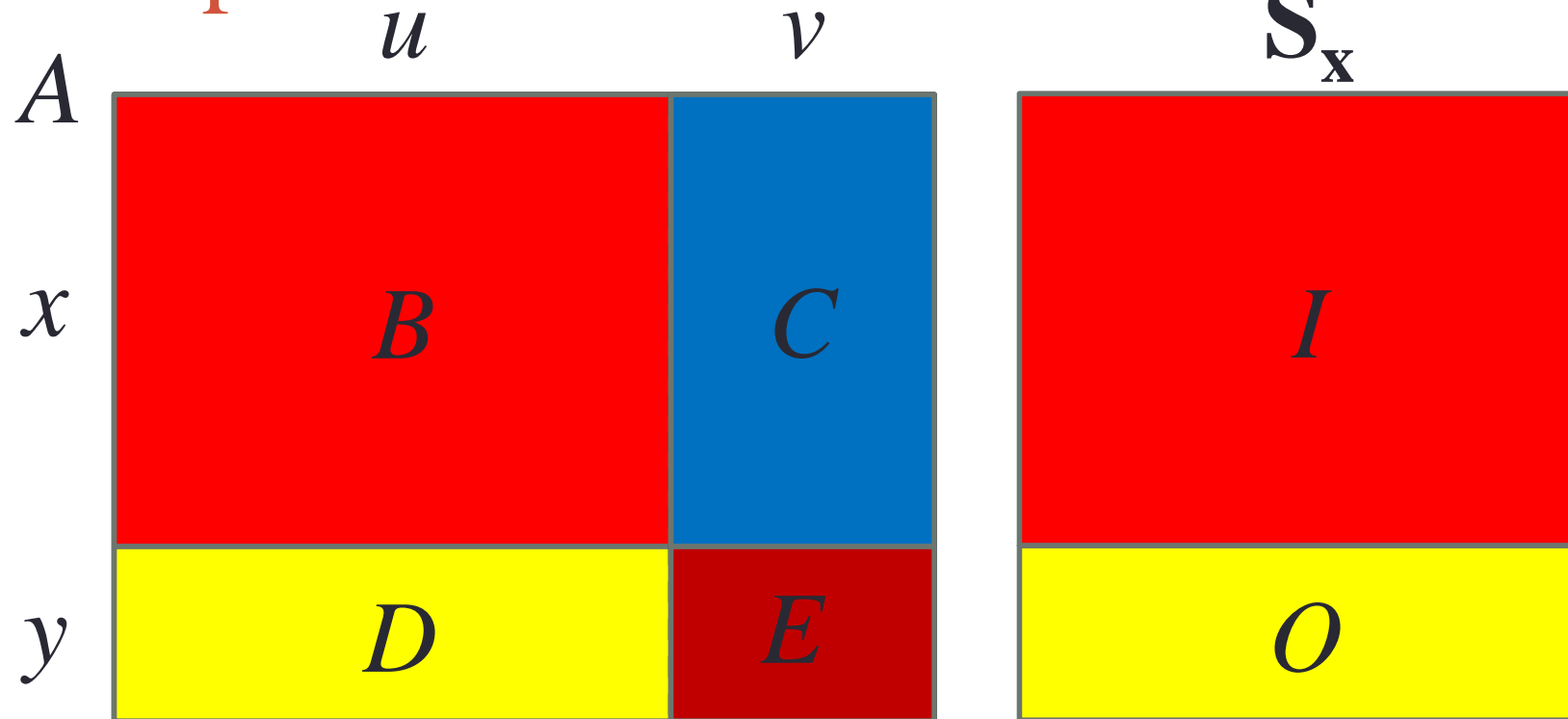
# Decomposition: selectors



$$A = \mathbf{S}_x \mathbf{B} \mathbf{T}_u + \mathbf{S}_x \mathbf{C} \mathbf{T}_v + \mathbf{S}_y \mathbf{D} \mathbf{T}_u + \mathbf{S}_y \mathbf{E} \mathbf{T}_v$$

with selectors  $\mathbf{S}_x \mathbf{S}_y \mathbf{T}_u \mathbf{T}_v$

# Decomposition: selectors

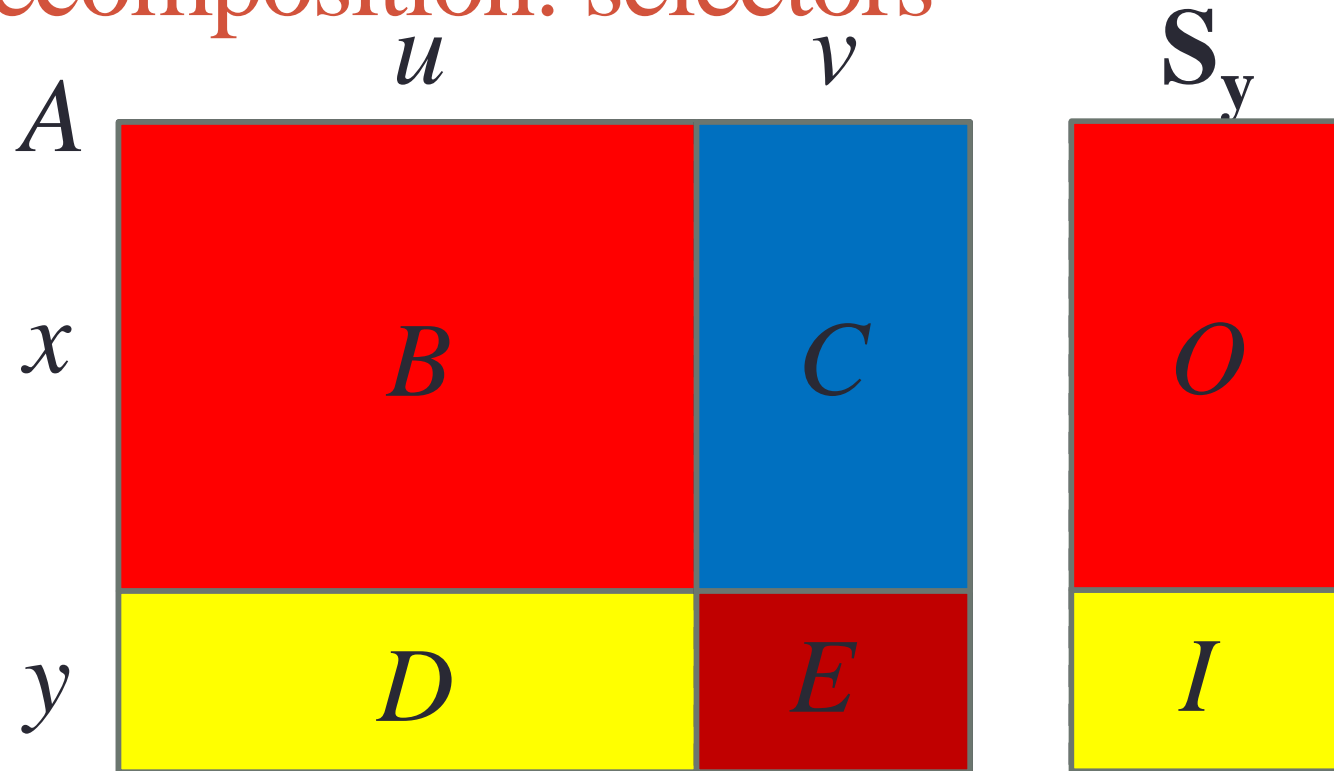


$$A = S_x B T_u + S_x C T_v + S_y D T_u + S_y E T_v$$

with selectors  $S_x S_y T_u T_v$



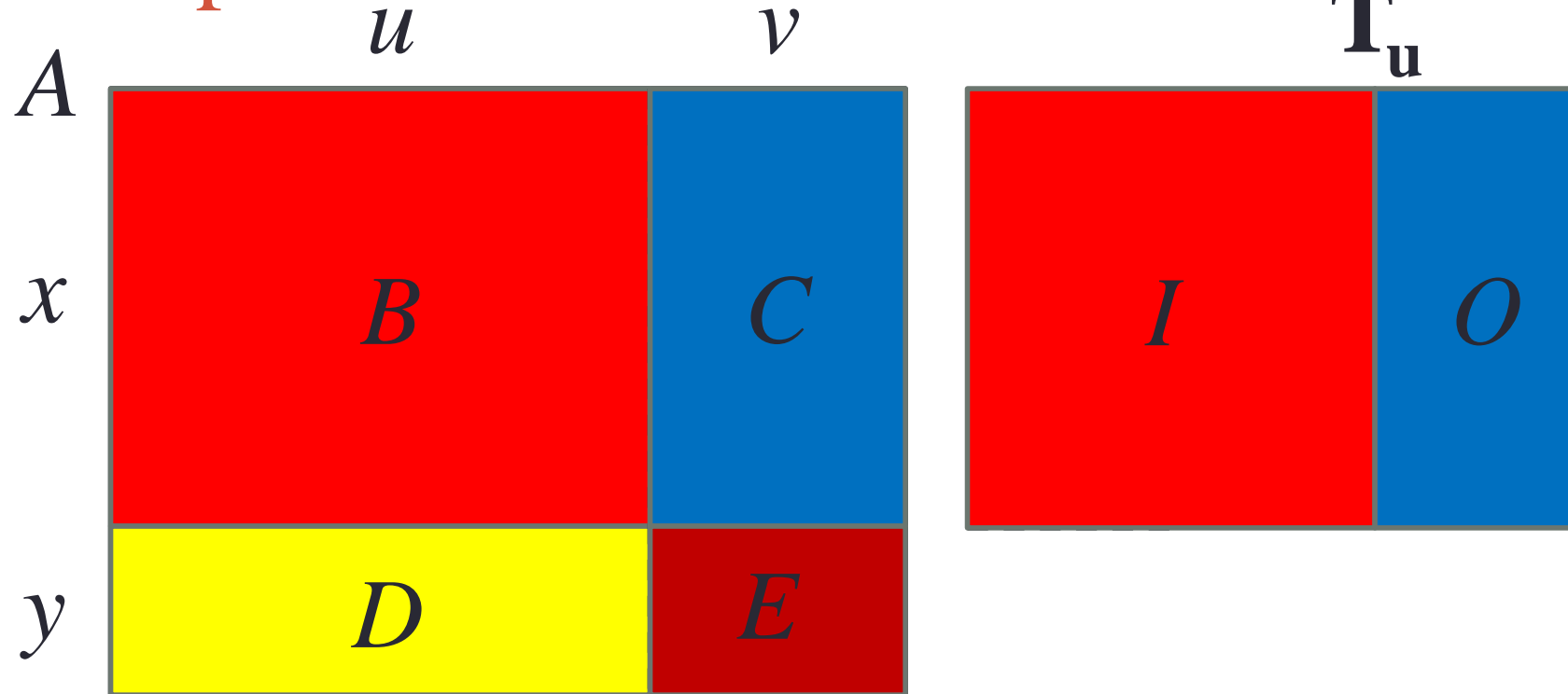
# Decomposition: selectors



$$A = S_x B T_u + S_x C T_v + S_y D T_u + S_y E T_v$$

with selectors  $S_x S_y T_u T_v$

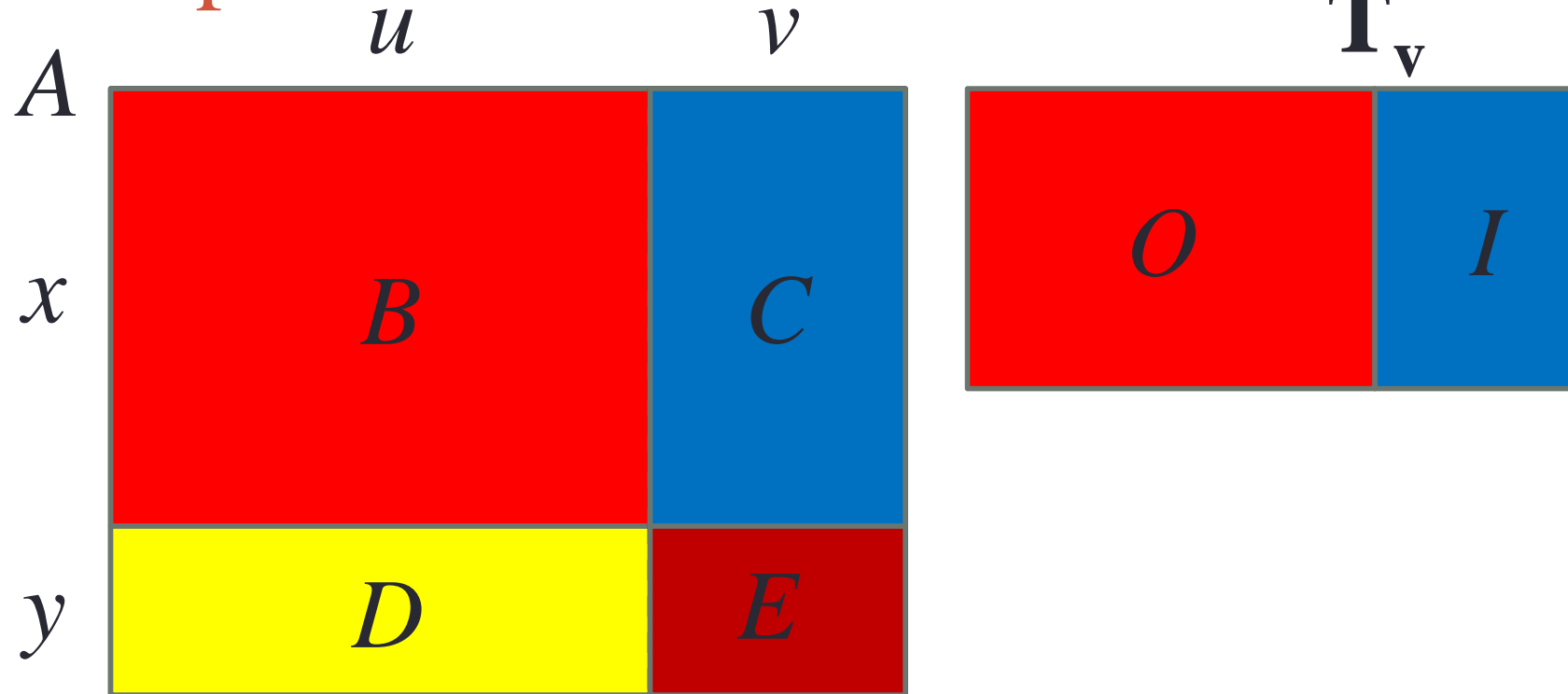
# Decomposition: selectors



$$A = S_x B T_u + S_x C T_v + S_y D T_u + S_y E T_v$$

with selectors  $S_x S_y T_u T_v$

# Decomposition: selectors



$$A = S_x B T_u + S_x C T_v + S_y D T_u + S_y E T_v$$

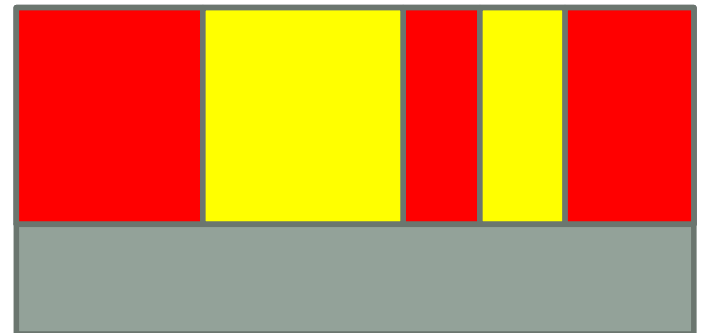
with selectors  $S_x S_y T_u T_v$

# Summary: selectors

- Relate a matrix to its sub-matrices.
- Matrices with only 0/1.
- Any row / column has at most a 1.

# Summary: selectors

- Relate a matrix to its sub-matrices.
- Matrices with only 0/1.
- Any row / column has at most a 1.
- Extends to **non-adjacent** case.
- Works for any **rectangular partition**.



# Provenance propagation

- We have
  - Matrices and operators over them – Algebra of matrices  $\mathcal{M}$
  - Annotations – Semiring of provenance polynomials  $\mathbb{N}[X]$
- Goals
  - Combine annotations in the same structure as the matrices
  - Operations should propagate data and annotations

# Provenance propagation

- We have
  - Matrices and operators over them – **Algebra of matrices**  $\mathcal{M}$
  - Annotations – **Semiring of provenance polynomials**  $\mathbb{N}[X]$
- Goals
  - Combine annotations in the same structure as the matrices
  - Operations should **propagate** data and annotations
- We do this in the space of tensor product  $\mathbb{N}[X] \otimes \mathcal{M}$ 
  - Matrices as vectors, provenance as scalars:  $p * A$
  - Satisfies all the laws of a  $\mathbb{N}[X]$ -**semialgebra**.

## Laws of a $N[X]$ -semialgebra ( $K$ -semialgebra)

$(K, +_K, \cdot_K, 0_K, 1_K)$  commutative semiring

$(K \otimes \mathcal{M}, +, \cdot, 0, 1)$  forms a ring (just like the matrices)

laws for scalar multiplication in a  $K$ -semialgebra

$$k * (A_1 + A_2) = k * A_1 + k * A_2$$

$$k * 0 = 0$$

$$(k_1 +_K k_2) * A = k_1 * A + k_2 * A$$

$$0_K * A = 0$$

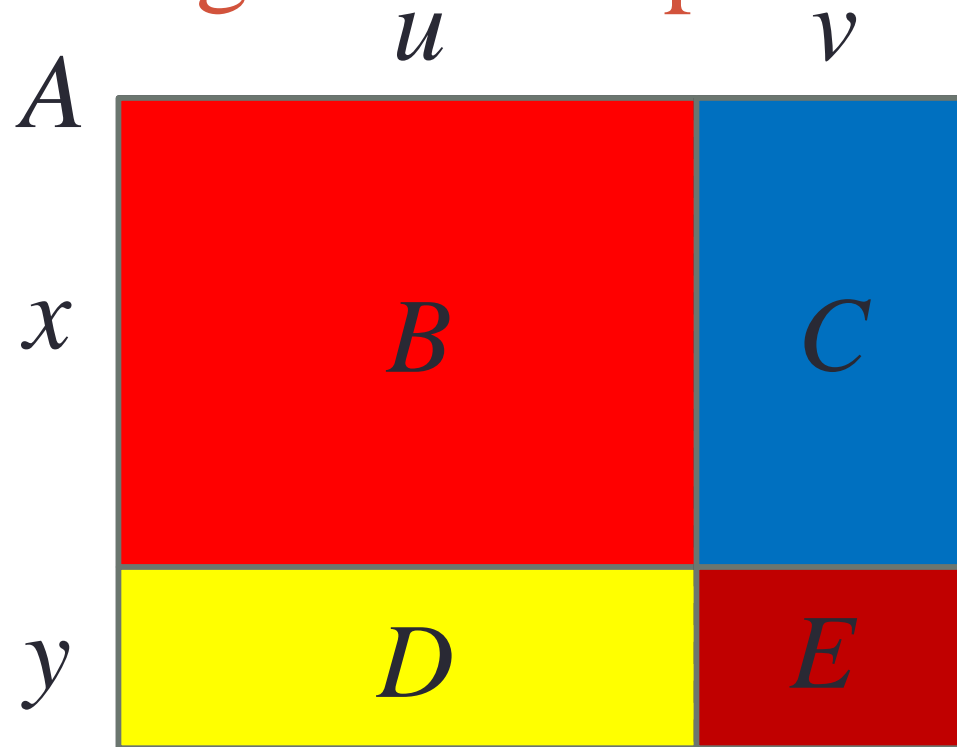
$$(k_1 \cdot_K k_2) * A = k_1 * (k_2 * A)$$

$$1_K * A = A$$

$$(k_1 * A_1)(k_2 * A_2) = (k_1 \cdot_K k_2) * (A_1 A_2)$$

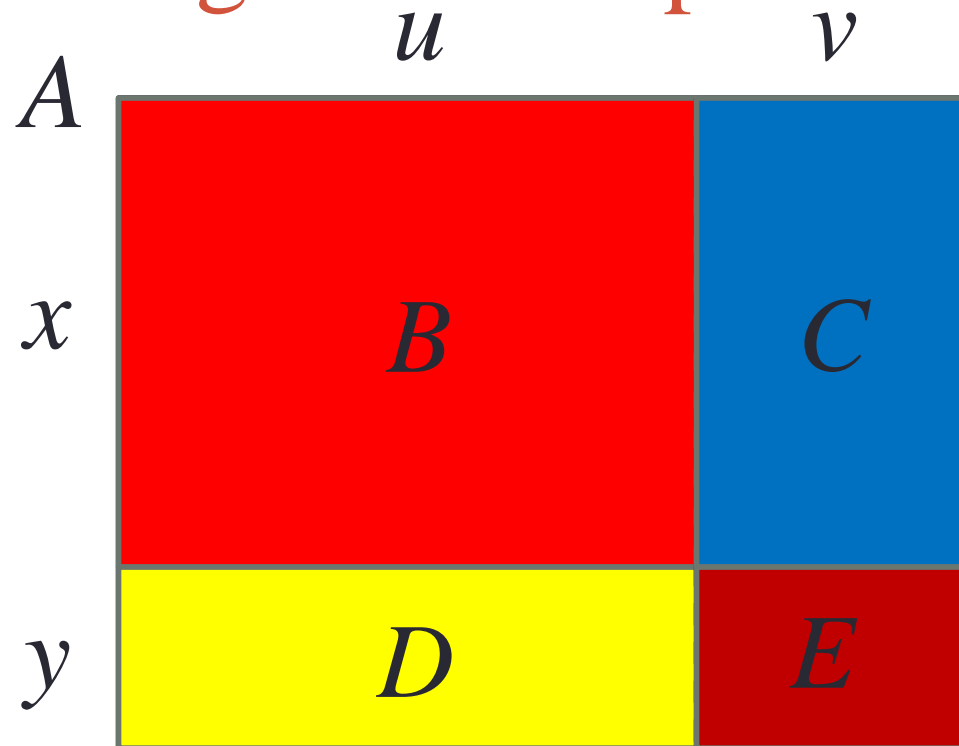


# Semialgebra example



$$A = \mathbf{S}_x B \mathbf{T}_u + \mathbf{S}_x C \mathbf{T}_v + \mathbf{S}_y D \mathbf{T}_u + \mathbf{S}_y E \mathbf{T}_v$$

# Semialgebra example: add annotation



$$\begin{aligned}
 A &= \mathbf{S}_x \text{ } xu^* B \mathbf{T}_u + \mathbf{S}_x \text{ } xv^* C \mathbf{T}_v \\
 &+ \mathbf{S}_y \text{ } yu^* D \mathbf{T}_u + \mathbf{S}_y \text{ } vy^* E \mathbf{T}_v
 \end{aligned}$$

## Semialgebra example: propagate annotation

$$\begin{aligned} A &= \mathbf{S}_x \textit{xu} * \mathbf{B} \mathbf{T}_u + \mathbf{S}_x \textit{xv} * \mathbf{C} \mathbf{T}_v \\ &+ \mathbf{S}_y \textit{yu} * \mathbf{D} \mathbf{T}_u + \mathbf{S}_y \textit{vy} * \mathbf{E} \mathbf{T}_v \end{aligned}$$

# Propagating annotation: transposition

$$A = \mathbf{S}_x \mathit{xu} * B \mathbf{T}_u + \mathbf{S}_x \mathit{xv} * C \mathbf{T}_v \\ + \mathbf{S}_y \mathit{yu} * D \mathbf{T}_u + \mathbf{S}_y \mathit{yv} * E \mathbf{T}_v$$

$$A^T = \mathbf{T}_u^T (\mathit{xu} * B^T) \mathbf{S}_x^T + \mathbf{T}_u^T (\mathit{yu} * D^T) \mathbf{S}_y^T \\ + \mathbf{T}_v^T (\mathit{xv} * C^T) \mathbf{S}_x^T + \mathbf{T}_v^T (\mathit{yv} * E^T) \mathbf{S}_y^T$$

# Propagating annotation: transposition

$$A = \mathbf{S}_x \text{ *xu* * } \mathbf{B} \mathbf{T}_u + \mathbf{S}_x \text{ *xv* * } \mathbf{C} \mathbf{T}_v \\ + \mathbf{S}_y \text{ *yu* * } \mathbf{D} \mathbf{T}_u + \mathbf{S}_y \text{ *vy* * } \mathbf{E} \mathbf{T}_v$$

$$A^T = \mathbf{T}_u^T (\text{*xu* * } \mathbf{B}^T) \mathbf{S}_x^T + \mathbf{T}_u^T (\text{*yu* * } \mathbf{D}^T) \mathbf{S}_y^T \\ + \mathbf{T}_v^T (\text{*xv* * } \mathbf{C}^T) \mathbf{S}_x^T + \mathbf{T}_v^T (\text{*vy* * } \mathbf{E}^T) \mathbf{S}_y^T$$

Transposition of a selector is still a selector  
 Still a sum of (selector  $\times$  matrix  $\times$  selector)

# Propagating annotation: multiplication

$$A = \mathbf{S}_x \mathit{xu} * B \mathbf{T}_u + \mathbf{S}_x \mathit{xv} * C \mathbf{T}_v \\ + \mathbf{S}_y \mathit{yu} * D \mathbf{T}_u + \mathbf{S}_y \mathit{vy} * E \mathbf{T}_v$$

$$A^T = \mathbf{T}_u^T (\mathit{xu} * B^T) \mathbf{S}_x^T + \mathbf{T}_u^T (\mathit{yu} * D^T) \mathbf{S}_y^T \\ + \mathbf{T}_v^T (\mathit{xv} * C^T) \mathbf{S}_x^T + \mathbf{T}_v^T (\mathit{yv} * E^T) \mathbf{S}_y^T$$

# Propagating annotation: multiplication

$$A = \mathbf{S}_x \mathit{xu} * B \mathbf{T}_u + \mathbf{S}_x \mathit{xv} * C \mathbf{T}_v \\ + \mathbf{S}_y \mathit{yu} * D \mathbf{T}_u + \mathbf{S}_y \mathit{vy} * E \mathbf{T}_v$$

$$A^T = \mathbf{T}_u^T (\mathit{xu} * B^T) \mathbf{S}_x^T + \mathbf{T}_u^T (\mathit{yu} * D^T) \mathbf{S}_y^T \\ + \mathbf{T}_v^T (\mathit{xv} * C^T) \mathbf{S}_x^T + \mathbf{T}_v^T (\mathit{yv} * E^T) \mathbf{S}_y^T$$

$$AA^T = \mathbf{S}_x (\mathit{x}^2 \mathit{u}^2 * BB^T + \mathit{x}^2 \mathit{v}^2 * CC^T) \mathbf{S}_x^T \\ + \mathbf{S}_x (\mathit{x} \mathit{y} \mathit{u}^2 * BD^T + \mathit{x} \mathit{y} \mathit{v}^2 * CE^T) \mathbf{S}_y^T \\ + \mathbf{S}_y (\mathit{x} \mathit{y} \mathit{u}^2 * DB^T + \mathit{x} \mathit{y} \mathit{v}^2 * EC^T) \mathbf{S}_x^T \\ + \mathbf{S}_y (\mathit{y}^2 \mathit{u}^2 * DD^T + \mathit{y}^2 \mathit{v}^2 * EE^T) \mathbf{S}_y^T$$

# Propagating annotation: multiplication

$$\begin{aligned}
 A &= \mathbf{S}_x \mathit{xu} * B \mathbf{T}_u + \mathbf{S}_x \mathit{xv} * C \mathbf{T}_v \\
 &+ \mathbf{S}_y \mathit{yu} * D \mathbf{T}_u + \mathbf{S}_y \mathit{vy} * E \mathbf{T}_v
 \end{aligned}
 \quad
 \begin{aligned}
 \mathbf{T}_u \mathbf{T}_u^T &= \mathbf{T}_v \mathbf{T}_v^T = \mathbf{I} \\
 \mathbf{T}_u \mathbf{T}_v^T &= \mathbf{T}_v \mathbf{T}_u^T = \mathbf{0}
 \end{aligned}$$

$$\begin{aligned}
 A^T &= \mathbf{T}_u^T (\mathit{xu} * B^T) \mathbf{S}_x^T + \mathbf{T}_u^T (\mathit{yu} * D^T) \mathbf{S}_y^T \\
 &+ \mathbf{T}_v^T (\mathit{xv} * C^T) \mathbf{S}_x^T + \mathbf{T}_v^T (\mathit{vy} * E^T) \mathbf{S}_y^T
 \end{aligned}$$

$$\begin{aligned}
 AA^T &= \mathbf{S}_x (x^2 u^2 * BB^T + x^2 v^2 * CC^T) \mathbf{S}_x^T \\
 &+ \mathbf{S}_x (xyu^2 * BD^T + xyv^2 * CE^T) \mathbf{S}_y^T \\
 &+ \mathbf{S}_y (xyu^2 * DB^T + xyv^2 * EC^T) \mathbf{S}_x^T \\
 &+ \mathbf{S}_y (y^2 u^2 * DD^T + y^2 v^2 * EE^T) \mathbf{S}_y^T
 \end{aligned}$$



## Semialgebra: delete propagation

$$\begin{aligned}
 AA^T &= \mathbf{S}_x(x^2u^2 * BB^T + x^2v^2 * CC^T)\mathbf{S}_x^T \\
 &+ \mathbf{S}_x(xy u^2 * BD^T + xy v^2 * CE^T)\mathbf{S}_y^T \\
 &+ \mathbf{S}_y(xy u^2 * DB^T + xy v^2 * EC^T)\mathbf{S}_x^T \\
 &+ \mathbf{S}_y(y^2u^2 * DD^T + y^2v^2 * EE^T)\mathbf{S}_y^T
 \end{aligned}$$

## Semialgebra: delete propagation

$$\begin{aligned}
 AA^T &= \mathbf{S}_x(x^2u^2 * BB^T + x^2v^2 * CC^T)\mathbf{S}_x^T \\
 &+ \mathbf{S}_x(xy u^2 * BD^T + xy v^2 * CE^T)\mathbf{S}_y^T \\
 &+ \mathbf{S}_y(xy u^2 * DB^T + xy v^2 * EC^T)\mathbf{S}_x^T \\
 &+ \mathbf{S}_y(y^2u^2 * DD^T + y^2v^2 * EE^T)\mathbf{S}_y^T
 \end{aligned}$$

**Deletion propagation:** set  $y = 0$

$$AA^T = x^2 * (u^2 * \mathbf{S}_x BB^T \mathbf{S}_x^T + v^2 * \mathbf{S}_x CC^T \mathbf{S}_x^T)$$

# Preliminary application: solving equations

- $(A + B)x = b$ ,  $A$  and  $B$  are square matrices
- $A$  is from source  $p$ ,  $B$  is from source  $q$

# Preliminary application: solving equations

- $(A + B)x = b$ ,  $A$  and  $B$  are square matrices
- $A$  is from source  $p$ ,  $B$  is from source  $q$
- **Jacobi method: iteratively compute**

$$u_{k+1} = (M^{-1}N)u_k + M^{-1}b \quad u_0 = \bar{0}$$

- $M = p * \text{diag}(A)$ ,  $N = p * (\text{diag}(A) - A) - q * B$

# Jacobi method: example

- Iteratively compute

$$u_{k+1} = (M^{-1}N)u_k + M^{-1}b \quad u_0 = \bar{0}$$

- $M = p * \text{diag}(A)$ ,  $N = p * (\text{diag}(A) - A) - q * B$

# Jacobi method: example

- Iteratively compute

$$u_{k+1} = (M^{-1}N)u_k + M^{-1}b \quad u_0 = \bar{0}$$

- $M = p * \text{diag}(A)$ ,  $N = p * (\text{diag}(A) - A) - q * B$

$$A = p * \begin{bmatrix} 2 & 1 \\ 1 & 2 \end{bmatrix}, B = q * \begin{bmatrix} 0 & 0 \\ -2 & 0 \end{bmatrix}, b = \begin{bmatrix} 1 \\ -1 \end{bmatrix}$$

$$u_1 = p * \begin{bmatrix} \frac{1}{2} \\ -\frac{1}{2} \end{bmatrix}, u_2 = p * \begin{bmatrix} \frac{1}{2} \\ -\frac{1}{2} \end{bmatrix} + p^3 * \begin{bmatrix} \frac{1}{4} \\ -\frac{1}{4} \end{bmatrix} + p^2 q * \begin{bmatrix} 0 \\ -\frac{1}{2} \end{bmatrix}, \dots$$

# Preliminary applications

- Solving systems of linear equations
- Also in the paper
  - Largest eigenvalue
  - PageRank

# Contributions

- First steps towards a **semantics-preserving** notion of **fine-grained provenance for linear algebra operators**
  - Key development: decomposition, tensor-product construction, and algebraic rules
- **Preliminary applications** in solving equations, computing largest eigenvalues, and PageRank.
- **Key benefit**
  - **Automatic propagation** of annotations through operators
  - **Ability to assign values** (e.g., 0 or 1) to the annotations and propagate the effects, e.g., for deletion or trust



# Related and future work

- Provenance Semirings / Semimodules
  - Green et al. PODS'07, Amsterdamer et al. PODS'11
- Array databases
  - SciDB, RasDaMan
  - Wu et al. SubZero, Peng and Diao SIGMOD'15
- Distributed machine learning / linear algebra programs
  - SystemML, Spark, MLbase, Cumulon, MADlib, GraphX, LINView, etc
- Future work
  - Support more linear algebra operators
  - Scalable implementation

Thank you!