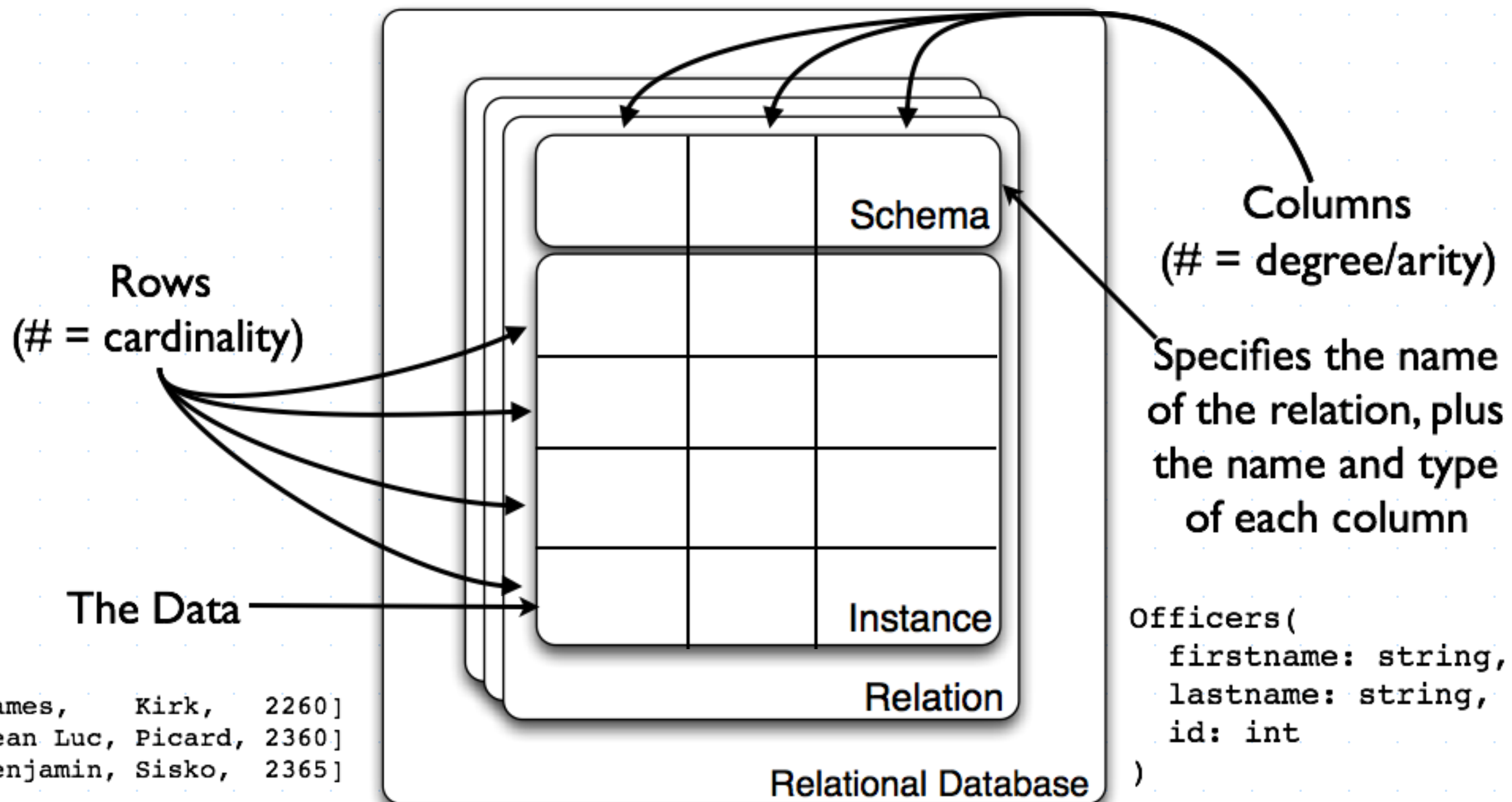


CSE 350

Advanced Data Structures

Lecture 2: Tidy Data



Example Data



Permits Economic & Neighborhood Development

This is a dataset of all permits issued from 1/2/2000 - present by the City of Buffalo. The Building Permit Office within the Department of Permit & Inspection Services is responsible for issuing permits, plan review and enforcement for compliance with the New York State Building Codes and all applicable ordinances in the City of Buffalo....

Last Updated
August 27, 2025

Data Provided By
Permit & Inspection Services

Featured Content Using this Data

<div>City of Buffalo Permit & Inspection Services</div> <div>External Content</div>	<div>🌐 Permit & Inspections Services</div> <div>Public</div> <div>April 15, 20251,927 Views</div>	<div>Permits & Applications</div> <div>External Content</div>
---	---	---

Explore: Examples/Permits.csv

1. What's the schema?

Permits, permit type, - - -

2. What's the arity?

37

3. What's the cardinality?

999

4. Are the ^{un}records of the table a set or a bag? Why?

Bag - No guarantee of uniqueness in a CSV file

Set - Conceptually unique way to reference records

4. If you're in charge of the CSV file, what could go wrong?

- Accidental Deletion
- Invalid Values
- Inconsistency

Goals for ...

Storage

No Redundancy

- Uniquely Identifiable data

Usage

Asking Questions

Functional Dependencies

Permit # Zip Code State Permit Type

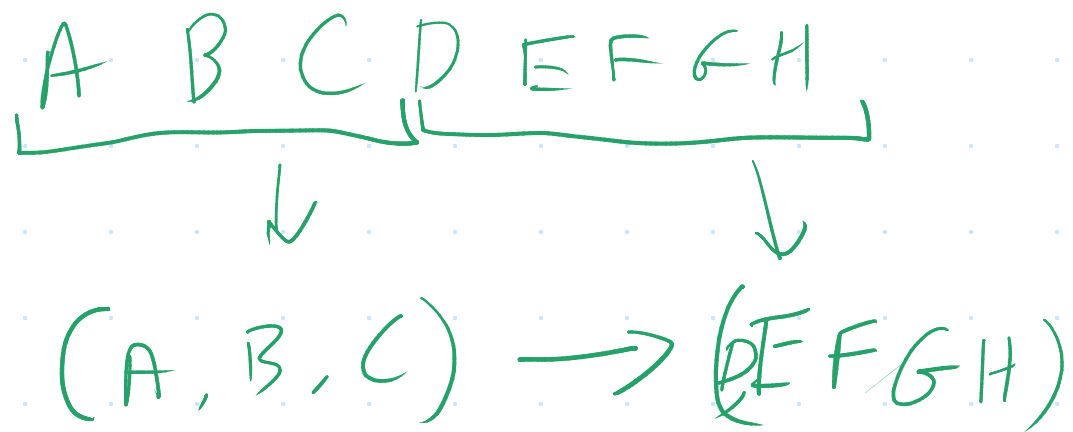
$\underline{A} \rightarrow B$

Permit #	Zip	State	Type
123	14228	N Y	Electrical
125	14228	N Y	Construction
127	10001	N Y	Electrical

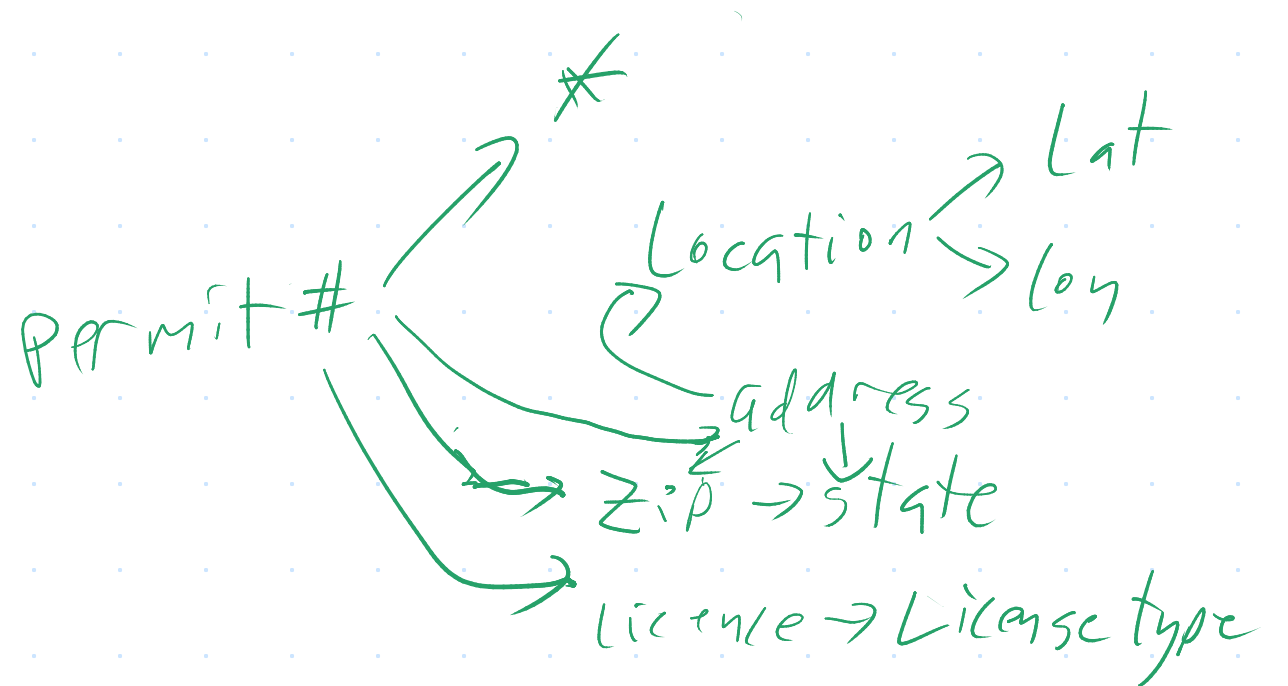
$\text{Permit} \# \rightarrow \text{Zip} \rightarrow \text{State}$
 $\text{Permit} \# \rightarrow \text{Zip} \rightarrow \text{Type}$
 $\text{Permit} \# \rightarrow \text{Zip}$

$\text{Zip} \rightarrow \text{State}$ (minimal)
 $\text{Permit} \rightarrow \text{Zip}$
 $\text{Permit} \rightarrow \text{State}$
 $\text{Permit} \rightarrow \text{Type}$

Superkeys

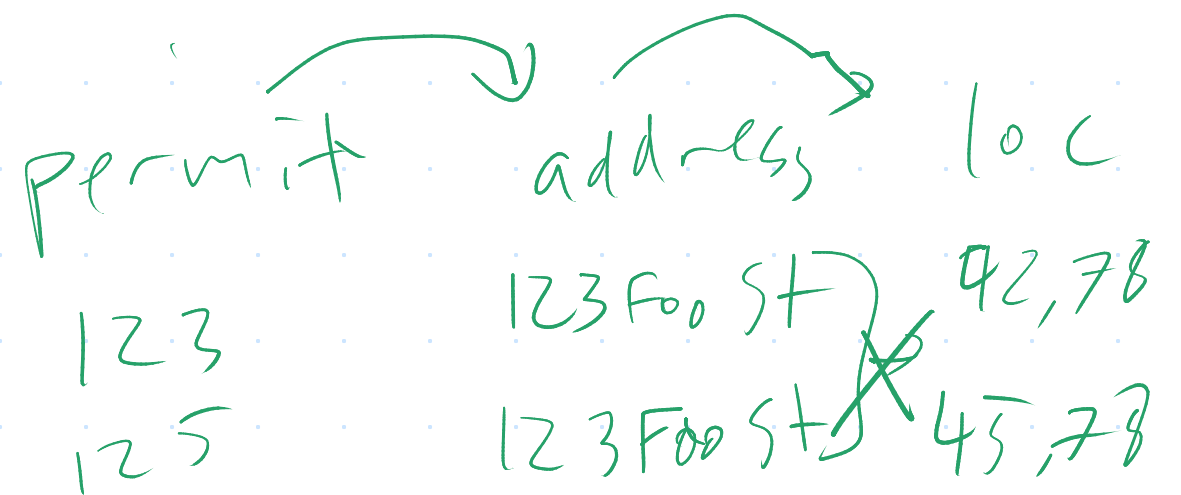


Candidate keys



permit #

Primary Keys



Rules of Tidy Data

1. One 'cell', one value
(is this broken anywhere?)

Location = Latitude & longitude

1 Normal Form

2. One key to rule them all
(is this broken anywhere?)

2 Normal Form

3. No delegation
(is this broken anywhere?)

~~permit → address → location~~

permit → address
permits.csv

2 Normal Form
~~address → loc~~
addresses.csv

Guidelines for Schema Design

1. What is the '~~thing~~^{entity}' that one record in the table is meant to model?

2. Can the '~~thing~~^{entity}' be broken down into smaller component things?

3. Does each '~~thing~~' have a name by which it can be uniquely identified?

4. Is there one canonical record where you would expect to find the '~~thing~~'?

Schema design for Permits.csv

Permits (PID, Type, Zip, State)

↳ Permits (PID, Type, Zip)

States (Zip, State)

The Raw Table ADT

Table

get_cursor()

insert(row)

Cursor

get() → row

next() → bool ← at end?

prev() → bool ←

update(row)

delete()

seek(position)

... but that's awkward. How about something easier?

Project (Map)

$\pi_{\text{city}}(\text{Permits}) \rightarrow \text{Table} \sim 1 \text{ row} \sim$
Per permit

```
output = []  
for row in Permits:  
    output.push(row["city"])  
return output
```

Select (Filter)

$\sigma_{\text{city} = \text{"Buffalo"}}(\text{Permits}) \rightarrow \text{Table} \sim \text{only}$
Buffalo permits

```
output = []  
for row in Permits:  
    if row["city"] == "Buffalo":  
        output.push(row)  
return output
```

Union

$A \cup B$

Join (Flatmap)

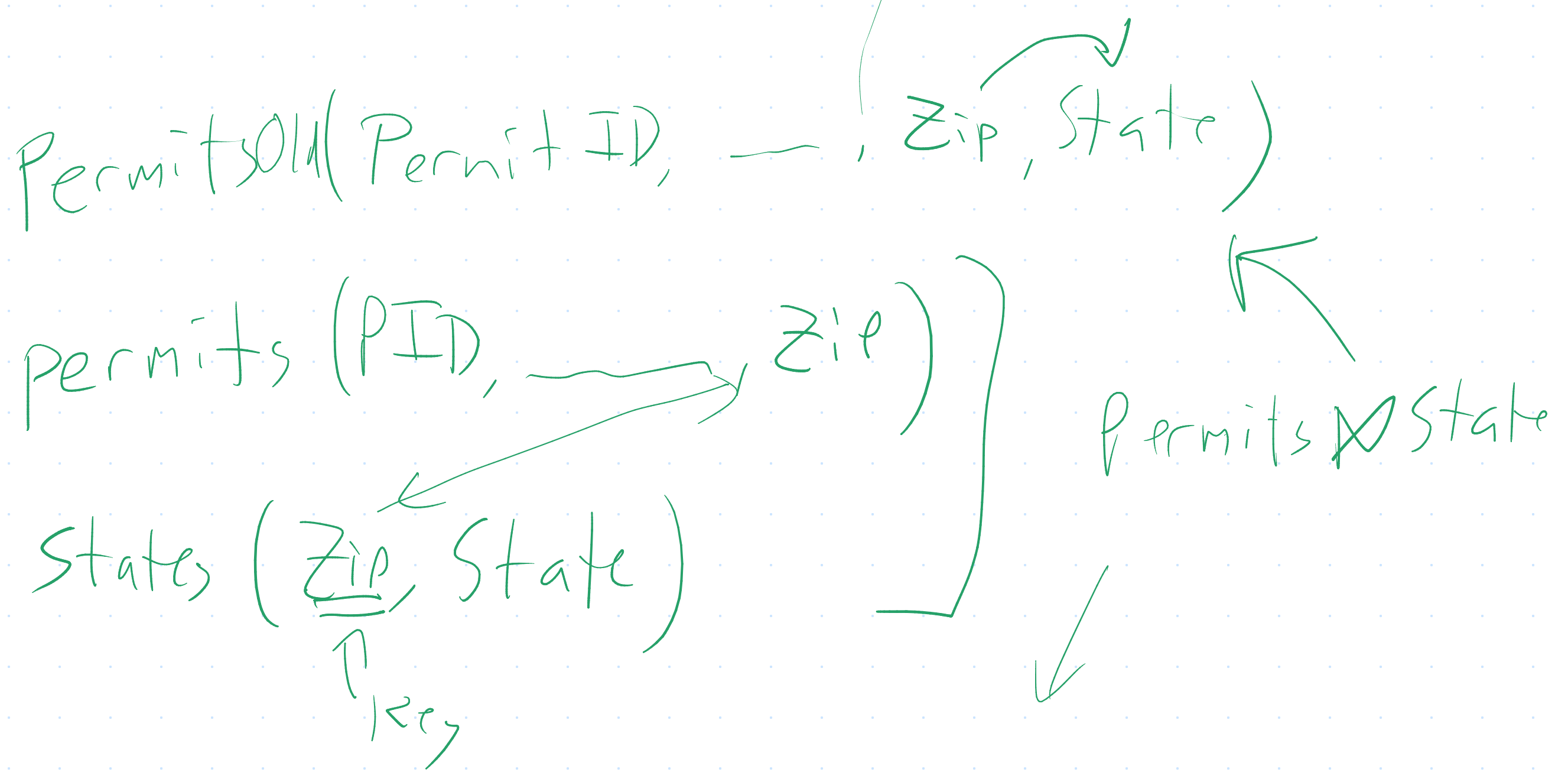
R	A	B	S	P	C
1	2		1	3	
1	3		2	4	
2	2		3	5	

output = []
 for r in R
 for s in S
 output.push(r + s)

R x S	A	B	S	P	C
1	2	1	3		
1	2	2	4		
1	2	3	5		
1	3	1	3		
1	3	2	4		
1	3	3	5		
2	2	1	3		
2	2	2	4		
2	2	3	5		

$$\sigma_{R.B=S.B}(R \times S) = R \bowtie_{R.B=S.B} S$$

"Equi Join"
 Implicit equality
 $R \bowtie S$ "Natural Join"



$$\begin{aligned}
 P \quad Q(PID, State) &= \pi_{\underline{PID}, \underline{State}} (Permits \bowtie States) \\
 &= \pi_{PID, State} (\pi_{PID, Zip} (Permits) \bowtie \pi_{Zip, State} (States))
 \end{aligned}$$

[Group-By] Aggregate (Fold)

$\{A, B, C, D, E\}$

init

accum

postprocess

$\rightarrow x$

$\rightarrow (value, x) \rightarrow x$

$\rightarrow x \rightarrow result$

$x = init()$

for r in R :

$x = accum(r, x)$

output(postprocess(x))

SUM
init

$= 0$

$accum(r, x) = r + x$

$postprocess(x) = x$

Average
init

$= \langle 0, 0 \rangle$

$accum(r, \langle x_c, x_s \rangle) = \langle x_c + r, x_s + r \rangle$

$postprocess(\langle x_c, x_s \rangle) = \frac{x_s}{x_c}$

$\sum_{\text{sum(Price)}} \text{Permits}$

$\sum_{\text{type, } \underline{\text{sum(Price)}}} \text{Permits}$

|
group by
attribute

|
aggregate

Others: Unique, Sort, Limit, Window

Why stick to these simple operators?

- Composable: Simple but expressive
- Easy to implement!
- Easy to implement efficiently
- Clear rules for when they can be rewritten (see 462)

... Even more friendly!

SELECT Applicant, Zipcode

FROM Permits

WHERE `Permit Type` = 'REPAIR'

$\Pi_{APP, ZIP} \left(\sigma_{\text{Permit Type} = \text{Repair}}(\text{permits}) \right)$

SELECT DISTINCT `Permit Type`

FROM Permits

ORDER BY

`Permit Type`

output must be a set

output must be a list

SELECT *
FROM Permits, Noise_Complaints
WHERE Permits.City = Noise_Complaints.City

~~X~~

$\sigma_{\text{Permits.city} = \text{Noise.city}}$ (permits X Noise)

SELECT COUNT(*)

FROM Permits, Noise_Complaints

WHERE Permits.City = Noise_Complaints.City

Σ (Permits.city | permits X Noise)
COUNT(*) = Noise.city

```
SELECT `Permit Type`, COUNT(*)
```

```
FROM Permits, Noise_Complaints
```

```
WHERE Permits.City = Noise_Complaints.City
```

```
GROUP BY 'Permit Type', Zip
```

$\sum_{\text{PermitType}, \text{COUNT}(*)} \left(\sum_{\text{Permits.city} = \text{Noise.city}} (\text{Permits} \times \text{Noise}) \right)$

More in a a few lectures...

Permits \bowtie [City] Noise_Complaints

```
for p in Permits
  for n in Noise
    if p.city = n.city
      output(p, n)
```

$O(|\text{Permits}| \cdot |\text{Noise}|)$

```
index = {}
for p in Permits
  index[p.city] += p
for n in Noise
  for p in index[n.city]
    output(p, n)
```

$O(|\text{Permits}| + |\text{Noise}| + |\text{output}|)$

$\Sigma[\text{Street}, \text{COUNT}(*)]$ (Permits)



Many of these

Big!

