Databases



COS 418/518: Distributed Systems
Lecture 19

Jialin Ding, Mike Freedman

Rest of the semester

- Lectures
 - Databases
 - Blockchains
 - Al systems
 - Reasoning about performance
- Assignments
 - Assignment 4 due Tuesday, November 25
 - Assignment 5 due Friday, December 12 (no late days)
- Final: Wednesday, December 17
 - Only covers material from the second half of the semester
 - More logistical details to come (will post to Ed)

Today

- Background on databases
- More on isolation
 - MVCC and snapshot isolation
- More on sharding
 - Data distribution strategies

Key-value stores vs. relational databases

- Key-value stores
 - Get(K)
 - Put(K, V)
 - Delete(K)

Key	Value
as1234	{"name": "Alice Smith", "major": "COS"}
bw5678	{"name": "Bob Williams", "major": "ECE"}

Key-value stores vs. relational databases

- Key-value stores
 - Get(K)
 - Put(K, V)
 - Delete(K)
- Relational databases
 - Tables (i.e., "relations") store records
 - Can select/filter records by different columns
 - Tables can be joined
 - Fixed schema

Key	Value
as1234	{"name": "Alice Smith", "major": "COS"}
bw5678	{"name": "Bob Williams", "major": "ECE"}

netid	first_name	last_name	major
as1234	Alice	Smith	COS
bw5678	Bob	Williams	ECE

major	enrollment
COS	400
ECE	300

Today

- Background on databases
- More on isolation
 - MVCC and snapshot isolation
- More on sharding
 - Data distribution strategies

Store multiple versions of each record

netid	first_name	last_name	major	version
as1234	Alice	Smith	COS	0
bw5678	Bob	Williams	ECE	0

Store multiple versions of each record

netid	first_name	last_name	major	version
as1234	Alice	Smith	COS	0
as1234	Alice	Smith	ECO	1
bw5678	Bob	Williams	ECE	0

- Store multiple versions of each record
 - Allows "time travel"
 - Enables snapshot isolation

netid	first_name	last_name	major	start_ts	end_ts
as1234	Alice	Smith	COS	Sep 2024	Oct 2025
as1234	Alice	Smith	ECO	Oct 2025	present
bw5678	Bob	Williams	ECE	Sep 2025	present

- Store multiple versions of each record
 - Allows "time travel"
 - Enables snapshot isolation

netid	first_name	last_name	major	start_xid	end_xid
as1234	Alice	Smith	COS	100	200
as1234	Alice	Smith	ECO	200	inf
bw5678	Bob	Williams	ECE	150	inf

Snapshot Isolation

- A transaction reads from a snapshot taken at start time
- No locks for reads!
- Conflicts may arise due to writes
 - If two transactions write the same record, the first transaction commits and the second transaction aborts (write-write conflict)
 - Does not check for read-write conflicts

T1:
$$R(A) - > 10$$
 $R(A) - > 20$

T2: W(A=20)

Snapshot Isolation

- A transaction reads from a snapshot taken at start time
- No locks for reads!
- Conflicts may arise due to writes
 - If two transactions write the same record, the first transaction commits and the second transaction aborts (write-write conflict)
 - Does not check for read-write conflicts

$$A = 100$$
 $A = 90$ $B = 100$ $A = 90$ $B = 110$ $B = 110$ $A = 90$ $B = 110$ $B = 110$ $A = 90$ $B = 110$ $B = 110$ $A = 90$ $B = 110$ $B = 110$ $A = 90$ $B = 110$ $B = 110$ $A = 90$ $B = 110$ $A = 90$ $B = 110$ $B = 110$ $A = 90$ $B = 110$ $B = 110$ $A = 90$ $B = 110$ $B = 110$ $A = 90$ $B = 110$ $B = 110$ $A = 90$ $B = 110$ $B = 110$ $A = 90$ $B = 110$ $B = 110$ $A = 90$ $B = 110$ $B = 110$ $A = 90$ $B = 110$ $B = 110$ $A = 90$ $B = 110$ $B = 110$ $A = 90$ $B = 110$ $B = 110$ $A = 90$ $B = 110$ $B = 110$ $A = 90$ $B = 110$ $B = 110$ $A = 90$ $B = 110$ $B = 110$ $A = 90$ $B = 110$ $B = 110$ $A = 90$ $B = 110$ $B = 110$ $A = 90$ $B = 110$ $B = 110$ $A = 90$ $B = 110$ $B = 110$ $A = 90$ $B = 110$ $A = 10$ $A = 90$ $A = 100$ $A = 10$

- Each transaction is committed with a monotonically increasing transaction ID (xid)
 - New records created by transaction will have xmin = xid
 - Records deleted by transaction will have xmax = xid

key	value	xmin	xmax
Α	10	0	inf

T1 (xid=1)

W(A, 20)

key	value	xmin	xmax
Α	10	0	1
Α	20	1	inf

T1 (xid=1)

W(A, 20)

- Each transaction is committed with a monotonically increasing transaction ID (xid)
 - New records created by transaction will have xmin = xid
 - Records deleted by transaction will have xmax = xid
- Snapshot is defined by three variables
 - xmin: oldest in-progress transaction ID at the time the snapshot was taken
 - xmax: next available transaction ID (not taken by committed or inprogress transactions)
 - xip: list of concurrent in-progress transaction IDs

nextXID = 1

key	value	xmin	xmax
Α	10	0	inf

T1

- Each transaction is committed with a monotonically increasing transaction ID (xid)
 - New records created by transaction will have xmin = xid
 - Records deleted by transaction will have xmax = xid
- Snapshot is defined by three variables
 - xmin: oldest in-progress transaction ID at the time the snapshot was taken
 - xmax: next available transaction ID (not taken by committed or in-progress transactions)
 - xip: list of concurrent in-progress transaction IDs
- Visibility check for a given record
 - Record must be created
 - xmin < snapshot xmin, OR
 - Snapshot xmin <= xmin < snapshot xmax AND xmin not in xip
 - Record must not be deleted
 - xmax >= snapshot xmax OR xmax in xip

nextXID = 1

key	value	xmin	xmax
Α	10	0	inf

T1

$$R(A) -> ?$$

Record must be created

- xmin < snapshot xmin, OR
- Snapshot xmin <= xmin < snapshot xmax AND xmin not in xip

Record must not be deleted

nextXID = 1

key	value	xmin	xmax
Α	10	0	inf

T1

$$R(A) -> 10$$

Record must be created

- xmin < snapshot xmin, OR
- Snapshot xmin <= xmin < snapshot xmax AND xmin not in xip

Record must not be deleted

nextXID = 2

key	value	xmin	xmax
Α	10	0	inf

$$xmin = 1$$

 $xmax = 1$

$$xip = {}$$

$$R(A) -> 10$$

Record must be created

- xmin < snapshot xmin, OR
- Snapshot xmin <= xmin < snapshot xmax AND xmin not in xip

Record must not be deleted

nextXID = 2

key	value	xmin	xmax
Α	10	0	1
Α	20	1	inf

T1 (xid=1)

xmin = 1 xmax = 1 xip = {}

R(A) -> 10 W(A, 20)

Record must be created

- xmin < snapshot xmin, OR
- Snapshot xmin <= xmin < snapshot xmax AND xmin not in xip

Record must not be deleted

n	extX		=	2
	IEXLAI	טו	_	_

key	value	xmin	xmax
Α	10	0	1
Α	20	1	inf

(' ')	
xmin = 1	xmin = 1
xmax = 1	xmax = 2
$xip = \{\}$	$xip = \{1\}$

R(A) -> 10 W(A, 20)

T1 (xid=1)

R(A) -> ?

T2

Record must be created

- xmin < snapshot xmin, OR
- Snapshot xmin <= xmin < snapshot xmax AND xmin not in xip

Record must not be deleted

nextX	\Box	=	2
$110\Lambda tM$		_	_

key	value	xmin	xmax
Α	10	0	1
Α	20	1	inf

,	
xmin = 1	xmin = 1
xmax = 1	xmax = 2
$xip = \{\}$	$xip = \{1\}$

T1 (xid=1)

R(A) -> 10

T2

Record must be created

- xmin < snapshot xmin, OR
- Snapshot xmin <= xmin < snapshot xmax AND xmin not in xip

Record must not be deleted

n	extX	ID	= 3
- 1 1	$I \cup X \cup X$		_ J

key	value	xmin	xmax
Α	10	0	1
Α	20	1	inf
В	30	2	inf

Record must be created

- xmin < snapshot xmin, OR
- Snapshot xmin <= xmin < snapshot xmax AND xmin not in xip

Record must not be deleted

T1 (xid=1)	T2 (xid=2)
xmin = 1 xmax = 1 xip = {}	xmin = 1 xmax = 2 xip = {1}
R(A) -> 10 W(A, 20)	
	R(A) -> 10
	W(B, 30)
	commit

n	extX	חו	= 3
ш	ICYLV	טו	

key	value	xmin	xmax
Α	10	0	1
Α	20	1	inf
В	30	2	inf

T1 (xid=1)	Т3
xmin = 1 xmax = 1	xmin = 1 xmax = 3
$xip = \{\}$	xip = {1}
R(A) -> 10 W(A, 20)	D(B)
	- (-)

Record must be created

- xmin < snapshot xmin, OR
- Snapshot xmin <= xmin < snapshot xmax AND xmin not in xip

Record must not be deleted

nextX	ID =	4
		_

key	value	xmin	xmax
Α	10	0	1
Α	20	1	inf
В	30	2	3

Record must be created

- xmin < snapshot xmin, OR
- Snapshot xmin <= xmin < snapshot xmax AND xmin not in xip

Record must not be deleted

T1 (xid=1) T3 (xid=3)

xmin = 1
xmax = 1
xip = {}

$$R(A) \rightarrow 10$$

W(A, 20)

D(B)

n	extX	=	4
	$10 \Lambda t \Lambda 1$		-

key	value	xmin	xmax
Α	10	0	1
Α	20	1	inf
В	30	2	3

Record	must	be o	created	
--------	------	------	---------	--

- xmin < snapshot xmin, OR
- Snapshot xmin <= xmin < snapshot xmax AND xmin not in xip

Record must not be deleted

T1 (xid=1) T3 (xid=3)

xmin = 1
xmax = 1
xip = {}

$$R(A) \rightarrow 10$$

W(A, 20)

 $R(A) \rightarrow 10$

D(B)
 $R(A) \rightarrow 10$

n	extXl	ID	=	1
	ロンスしへに			4

W(A, 40)

key	value	xmin	xmax
А	10	0	1
Α	20	1	inf
В	30	2	3

Record must be created

- xmin < snapshot xmin, OR
- Snapshot xmin <= xmin < snapshot xmax AND xmin not in xip

Record must not be deleted

T1 (xid=1)	T3 (xid=3)
xmin = 1	xmin = 1
xmax = 1	xmax = 3
xip = {}	xip = {1}
R(A) -> 10	D(B)
W(A, 20)	R(A) -> 10

Snapshot Isolation vs. Serializability

- Serializability is a stronger isolation level
 - Snapshot isolation allows schedules that are not serializable (next slide)
- Snapshot isolation is more scalable
 - Reads and writes do not block each other, unlike 2PL

Write skew

- Occurs when transactions read overlapping sets of data but write disjoint sets
 - Snapshot isolation only checks for write-write conflicts on the same record

key	value
Α	1
В	0

T1 (turn all 1s to 0s)	T2 (turn all 0s to 1s)
R(A)	R(A)
R(B)	R(B)
W(A, 0)	W(B, 1)

Write skew

- Occurs when transactions read overlapping sets of data but write disjoint sets
 - Snapshot isolation only checks for write-write conflicts on the same record

key	value
Α	1
В	0

T1 (turn all 1s to 0s)	T2 (turn all 0s to 1s)
R(A)	R(A)
R(B)	R(B)
W(A, 0)	W(B, 1)

key	value
Α	0
В	1

Not serializable!

Isolation vs. consistency

- Isolation is about whether concurrent transactions interfere with each other
 - Isolation models: serializable, snapshot isolation
- Consistency is about whether nodes see the same state
 - Consistency models: linearizable, causal+, eventual

Today

- Background on databases
- More on isolation
 - MVCC and snapshot isolation
- More on sharding
 - Data distribution strategies

- Words that mean similar things
 - Sharding
 - Partitioning
 - Careful: could refer to data partitioning or network partitioning
 - Distribution
 - Careful: data can be distributed without being sharded

	Node 1	Node 2	Node 3
Range	Netid in [aa0000, hz9999]	Netid in [ia0000, sz0000]	Netid in [ta0000, zz9999]
Round robin			
Hash			
All/broadcast			

	Node 1	Node 2	Node 3
Range	Netid in [aa0000, hz9999]	Netid in [ia0000, sz0000]	Netid in [ta0000, zz9999]
Round robin	Row 1, row 4, row 7,	Row 2, row 5, row 8,	Row 3, row 6, row 9,
Hash			
All/broadcast			

	Node 1	Node 2	Node 3
Range	Netid in [aa0000, hz9999]	Netid in [ia0000, sz0000]	Netid in [ta0000, zz9999]
Round robin	Row 1, row 4, row 7,	Row 2, row 5, row 8,	Row 3, row 6, row 9,
Hash	Hash(netid) mod 3 = 0	Hash(netid) mod 3 = 1	Hash(netid) mod 3 = 2
All/broadcast			

	Node 1	Node 2	Node 3
Range	Netid in [aa0000, hz9999]	Netid in [ia0000, sz0000]	Netid in [ta0000, zz9999]
Round robin	Row 1, row 4, row 7,	Row 2, row 5, row 8,	Row 3, row 6, row 9,
Hash	Hash(netid) mod 3 = 0	Hash(netid) mod 3 = 1	Hash(netid) mod 3 = 2
All/broadcast	All	All	All

	Range	Round Robin	Hash	All/broadcast
Need to select a column?	Yes	No	Yes	No
Load balancing	Hard	Easy	Easy	N/A

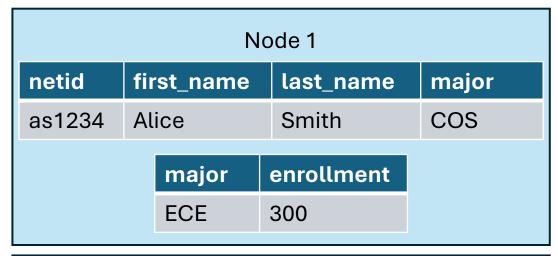
	Range	Round Robin	Hash	All/broadcast
Need to select a column?	Yes	No	Yes	No
Load balancing	Hard	Easy	Easy	N/A
Impact on joins	Can be useful but makes load balancing worse	Useless	Good for large tables	Good for small tables

netid	first_name	last_name	major
as1234	Alice	Smith	COS
bw5678	Bob	Williams	ECE

major	enrollment
ECE	300
COS	400

Round robin distribution

Joined records are not co-located on the same node -> network communication overhead



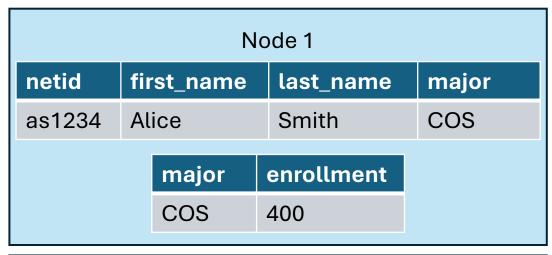
Node 2					
netid	fir	st_name		last_name	major
bw5678	Bob			Williams	ECE
		major	•	enrollment	
		COS	4	400	

netid	first_name	last_name	major
as1234	Alice	Smith	COS
bw5678	Bob	Williams	ECE

major	enrollment
ECE	300
COS	400

Hash distribution (on "major" column for both tables)

Joined records are co-located on the same node -> No network communication overhead

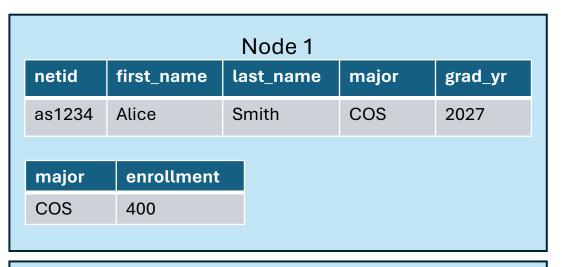


Node 2				
first_name		netid first_name last_name		major
Bob			Williams	ECE
	major		enrollment	
	ECE	,	300	
		first_name Bob major	first_name Bob major	first_name last_name Bob Williams major enrollment

netid	first_name	last_name	major	grad_yr
as1234	Alice	Smith	COS	2027
bw5678	Bob	Williams	ECE	2026

major	enrollment
ECE	300
COS	400

grad_yr	enrollment
2026	1500
2027	1600



		Node 2		
netid	first_name	last_name	major	grad_yr
bw5678	Bob	Williams	ECE	2026
major	enrollment			
ECE	300			

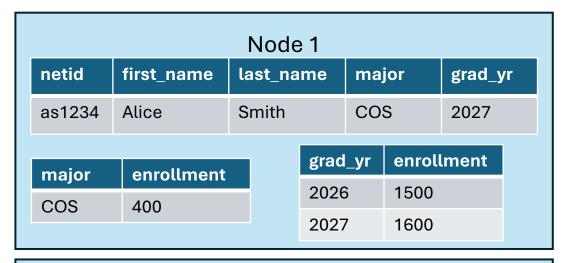
netid	first_name	last_name	major	grad_yr
as1234	Alice	Smith	COS	2027
bw5678	Bob	Williams	ECE	2026

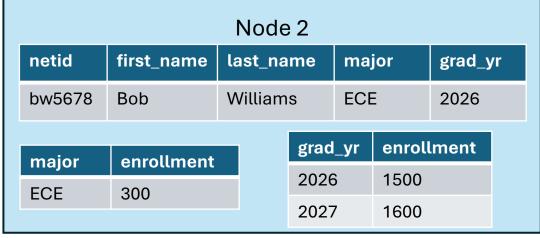
major	enrollment
ECE	300
COS	400

grad_yr	enrollment
2026	1500
2027	1600

All/broadcast distribution for small tables

Data is duplicated but network communication is minimized





Today

- Background on databases
- More on isolation
 - MVCC and snapshot isolation
- More on sharding
 - Data distribution strategies