# Scaling up
# Correlation Clustering
# through Parallelism and Concurrency Control



Xinghao Pan     Dimitris Papailiopoulos     Benjaminn Recht     Kannan Ramchandran     Michael I. Jordan
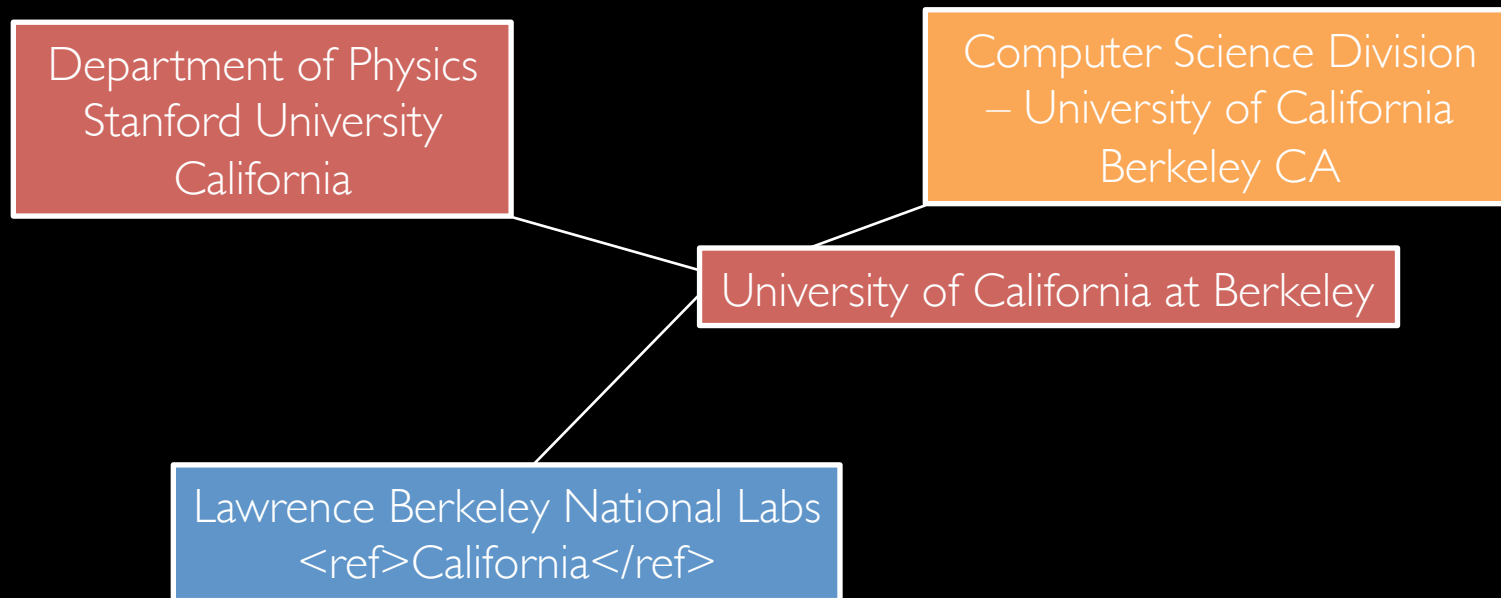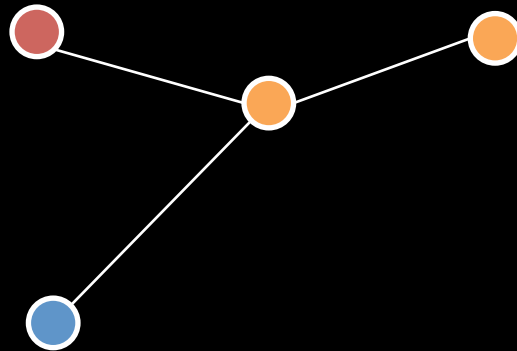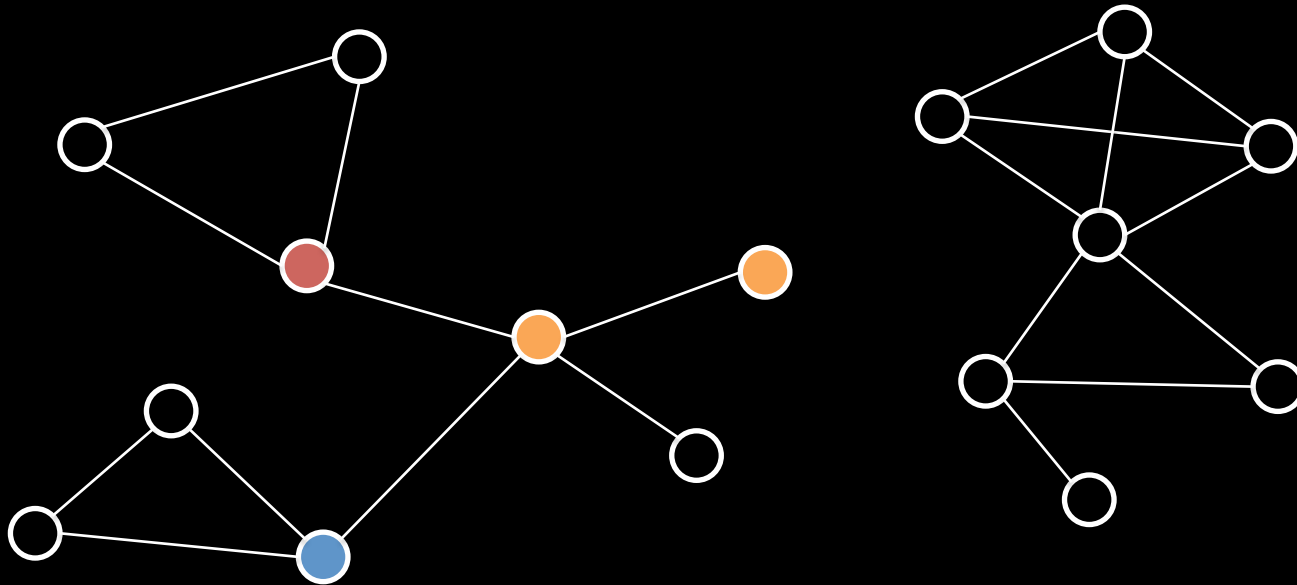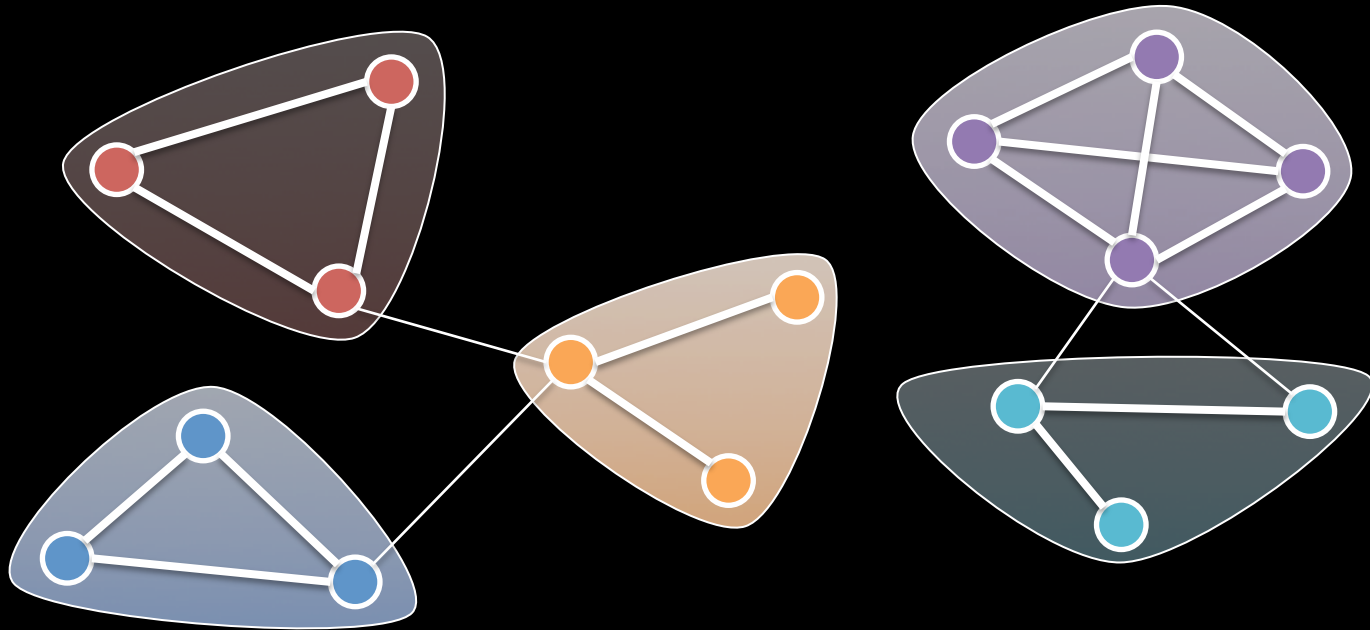
# Correlation Clustering for Deduplication

# Correlation Clustering for Deduplication

# Correlation Clustering for Deduplication

# Correlation Clustering for Deduplication

# Serial Correlation Clustering

Nir Ailon, Moses Charikar, and Alantha Newman.
Aggregating inconsistent information: ranking and clustering.
*Journal of the ACM (JACM)*, 55(5):23, 2008.

# Serial Correlation Clustering

Serially process vertices

Nir Ailon, Moses Charikar, and Alantha Newman.
Aggregating inconsistent information: ranking and clustering.
*Journal of the ACM (JACM)*, 55(5):23, 2008.

# Serial Correlation Clustering

Serially process vertices

Nir Ailon, Moses Charikar, and Alantha Newman.
Aggregating inconsistent information: ranking and clustering.
*Journal of the ACM (JACM)*, 55(5):23, 2008.

# Serial Correlation Clustering

Serially process vertices

Nir Ailon, Moses Charikar, and Alantha Newman.
Aggregating inconsistent information: ranking and clustering.
*Journal of the ACM (JACM)*, 55(5):23, 2008.

# Serial Correlation Clustering

Serially process vertices



Approximation 3 OPT (in expectation)

Nir Ailon, Moses Charikar, and Alantha Newman.
Aggregating inconsistent information: ranking and clustering.
*Journal of the ACM (JACM)*, 55(5):23, 2008.

# Serial Correlation Clustering

Serially process vertices

Objective: Parallelize Correlation Clustering
with strong guarantees and high concurrency.

Idea:  Operations on vertices ⇔ Database transaction.
Apply concepts of concurrency control from databases.

Approximation 3 OPT (in expectation)

Nir Ailon, Moses Charikar, and Alantha Newman.
Aggregating inconsistent information: ranking and clustering.
*Journal of the ACM (JACM)*, 55(5):23, 2008.

# Machine Learning + Concurrency Control
(Xinghao Pan et al.)

Serial ML algorithm

# Machine Learning + Concurrency Control

(Xinghao Pan et al.)

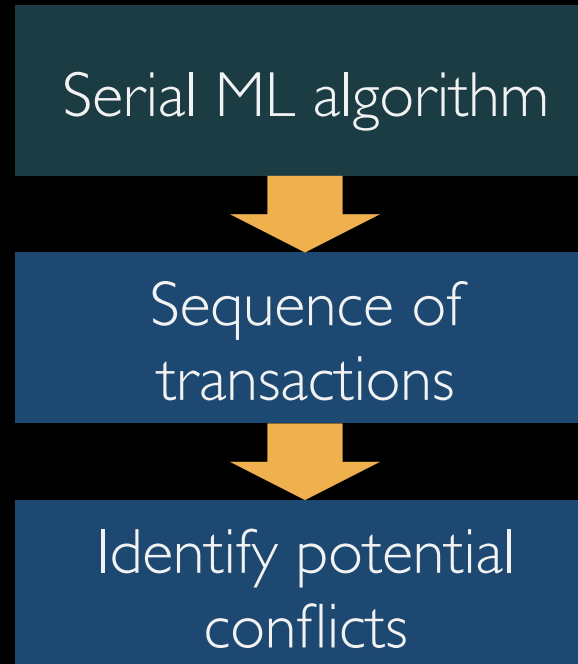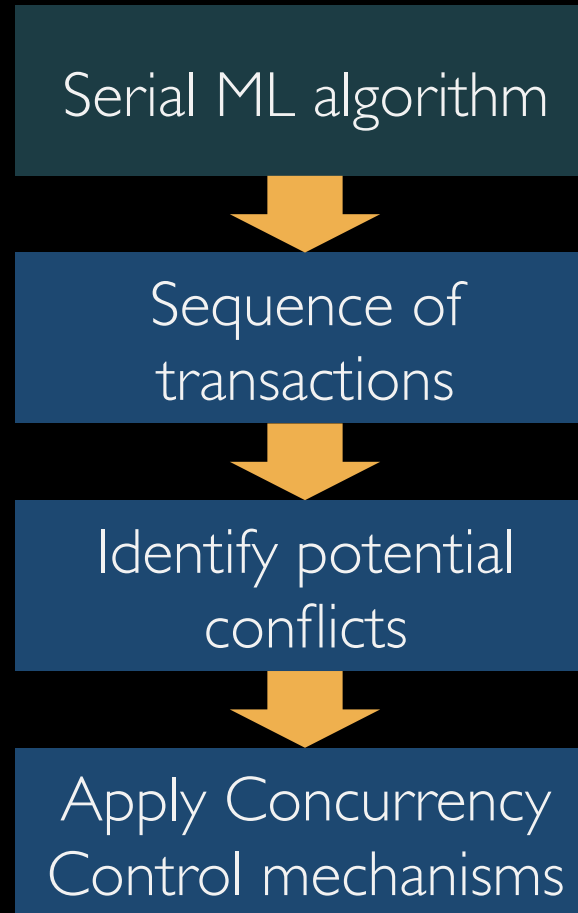Serial ML algorithm

↓

Sequence of transactions

# Machine Learning + Concurrency Control

(Xinghao Pan et al.)

Serial ML algorithm

⬇

Sequence of transactions

⬇

Identify potential conflicts

# Machine Learning + Concurrency Control
(Xinghao Pan et al.)

```
┌─────────────────────────┐
│   Serial ML algorithm   │
└─────────────────────────┘
            ↓
┌─────────────────────────┐
│      Sequence of        │
│      transactions       │
└─────────────────────────┘
            ↓
┌─────────────────────────┐
│   Identify potential    │
│       conflicts         │
└─────────────────────────┘
            ↓
┌─────────────────────────┐
│   Apply Concurrency     │
│  Control mechanisms     │
└─────────────────────────┘
```

# Machine Learning + Concurrency Control
(Xinghao Pan et al.)

# Properties of C4
## (Concurrency Control Correlation Clustering)

Theorem: C4 is **correct**.

C4 <u>preserves same guarantees</u> as serial algorithm (3 OPT).

# Properties of C4
## (Concurrency Control Correlation Clustering)

**Correctness**

**Theorem:** C4 is **correct**.

C4 <u>preserves same guarantees</u> as serial algorithm (3 OPT).

**Concurrency**

**Theorem:** C4 has **small overheads.**

= almost **linear speedup**

Expected #blocked transactions $< 2\tau\, |E|\, /\, |V|$.

$\tau \equiv$ diff in parallel cpu's progress

# Properties of C4
## (Concurrency Control Correlation Clustering)
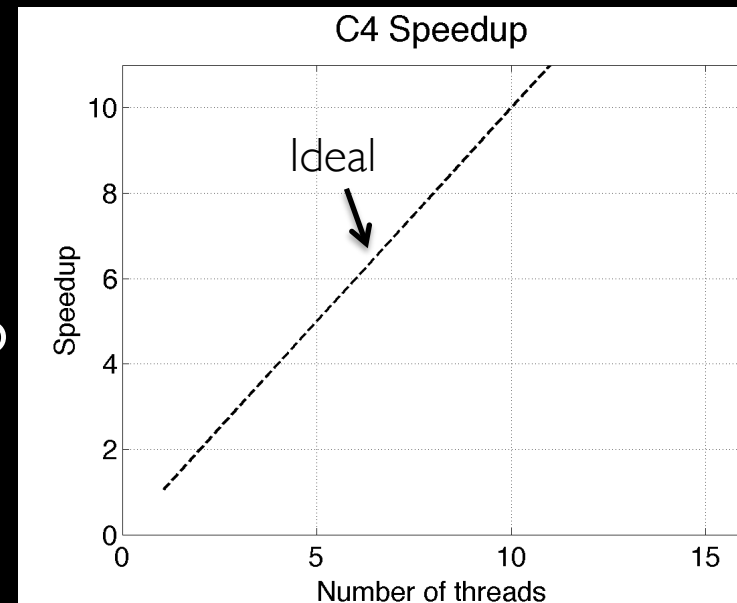
**Correctness**

Theorem: C4 is **correct**.

C4 <u>preserves same guarantees</u> as serial algorithm (3 OPT).

**Concurrency**

Theorem: C4 has **small overheads.**
       = almost **linear speedup**

Expected #blocked transactions $< 2\tau\, |E| / |V|$.

$\tau \equiv$ diff in parallel cpu's progress



C4 Speedup

Ideal

Speedup

Number of threads

# Properties of C4
## (Concurrency Control Correlation Clustering)

**Correctness**

Theorem: C4 is **correct**.

C4 <u>preserves same guarantees</u> as serial algorithm (3 OPT).

**Concurrency**

Theorem: C4 has **small overheads.**
        = almost **linear speedup**

Expected #blocked transactions < 2$\tau$ $|E|$ / $|V|$.

$\tau$ ≡ diff in parallel cpu's progress

### C4 Speedup

Ideal

10x
Speed up

- - - - Ideal
○····· IT-2004
✕— Webbase-2001
☐— Erdos-Renyi

Speedup (y-axis): 0, 2, 4, 6, 8, 10

Number of threads (x-axis): 0, 5, 10, 15