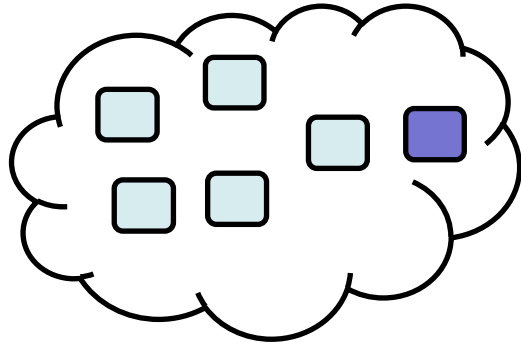


Update Propagation in a Streaming Warehouse

Theodore Johnson
Vladislav Shkapenyuk

SSDBM 2011

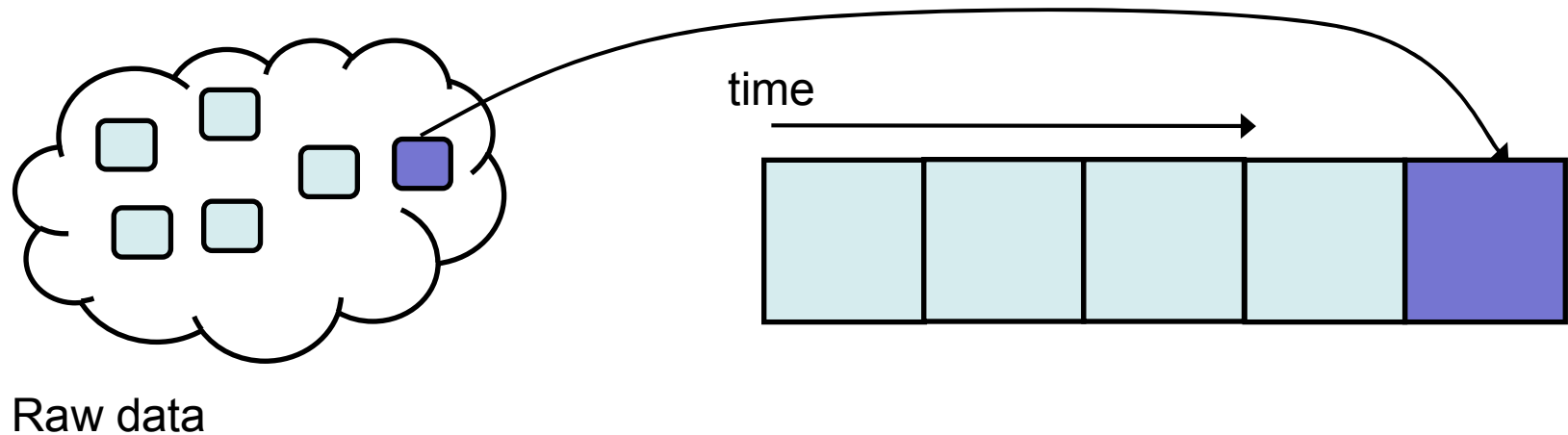
Managing a Stream Warehouse



Raw data

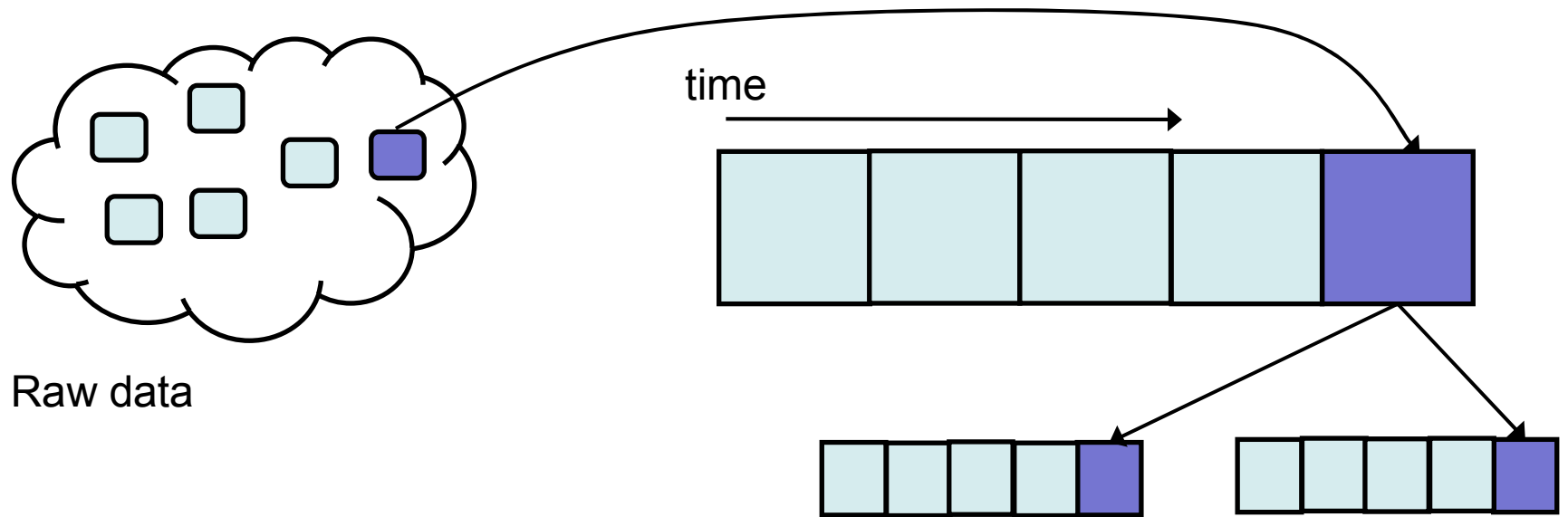
- Continually arriving data

Managing a Stream Warehouse



- Continually arriving data
- Is loaded into temporally partitioned base tables

Managing a Stream Warehouse

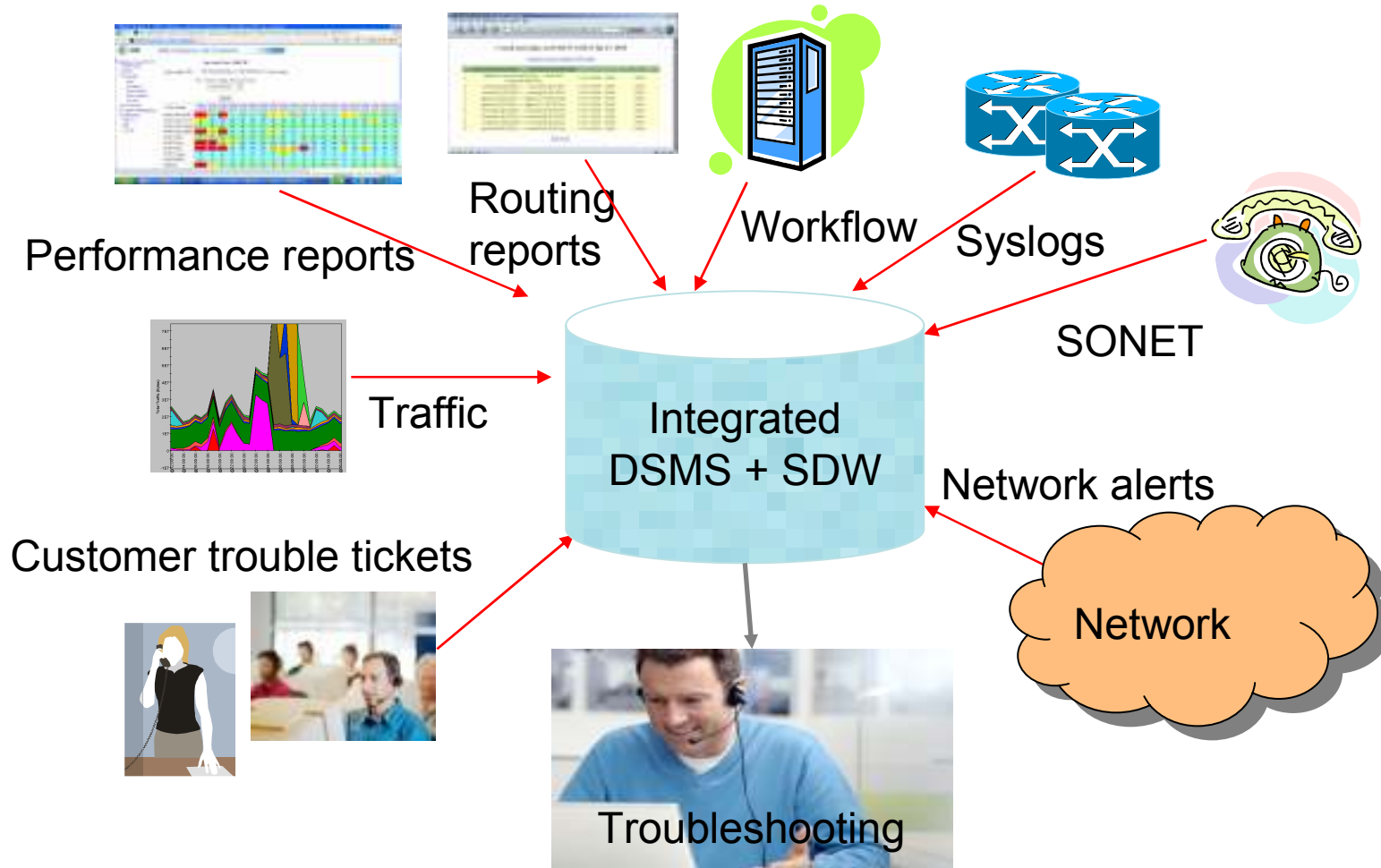


- Continually arriving data
- Is loaded into temporally partitioned base tables
- Updates propagate to higher level data products.

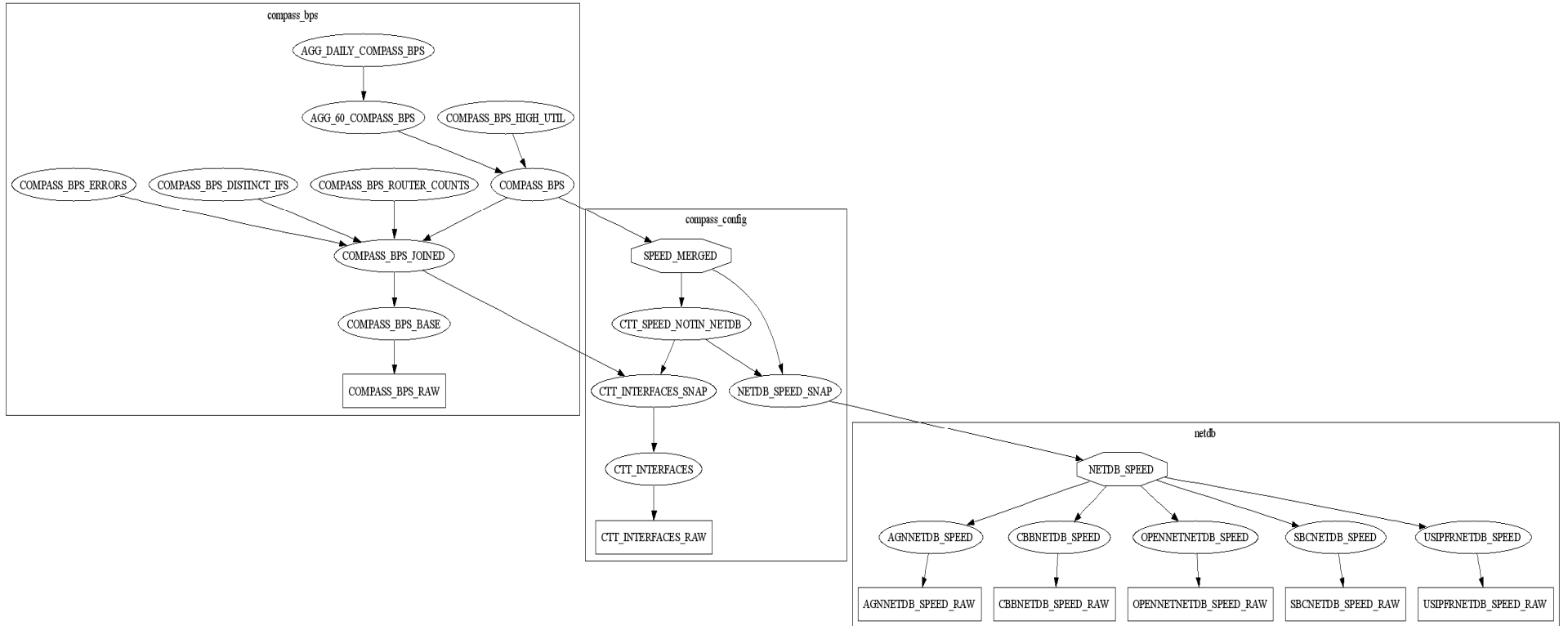
DataDepot

- DataDepot is a stream warehouse generator
 - On top of the Daytona database.
- Applications within AT&T
 - Darkstar warehouse: Network performance, configuration, and alert data.
 - Other petabyte-sized applications.
- Materialized view maintenance via update propagation through partitions.
- Near-real time: Propagate updates as soon as possible
 - Don't wait for global updates to materialized views.
- Multiversion Concurrency Control

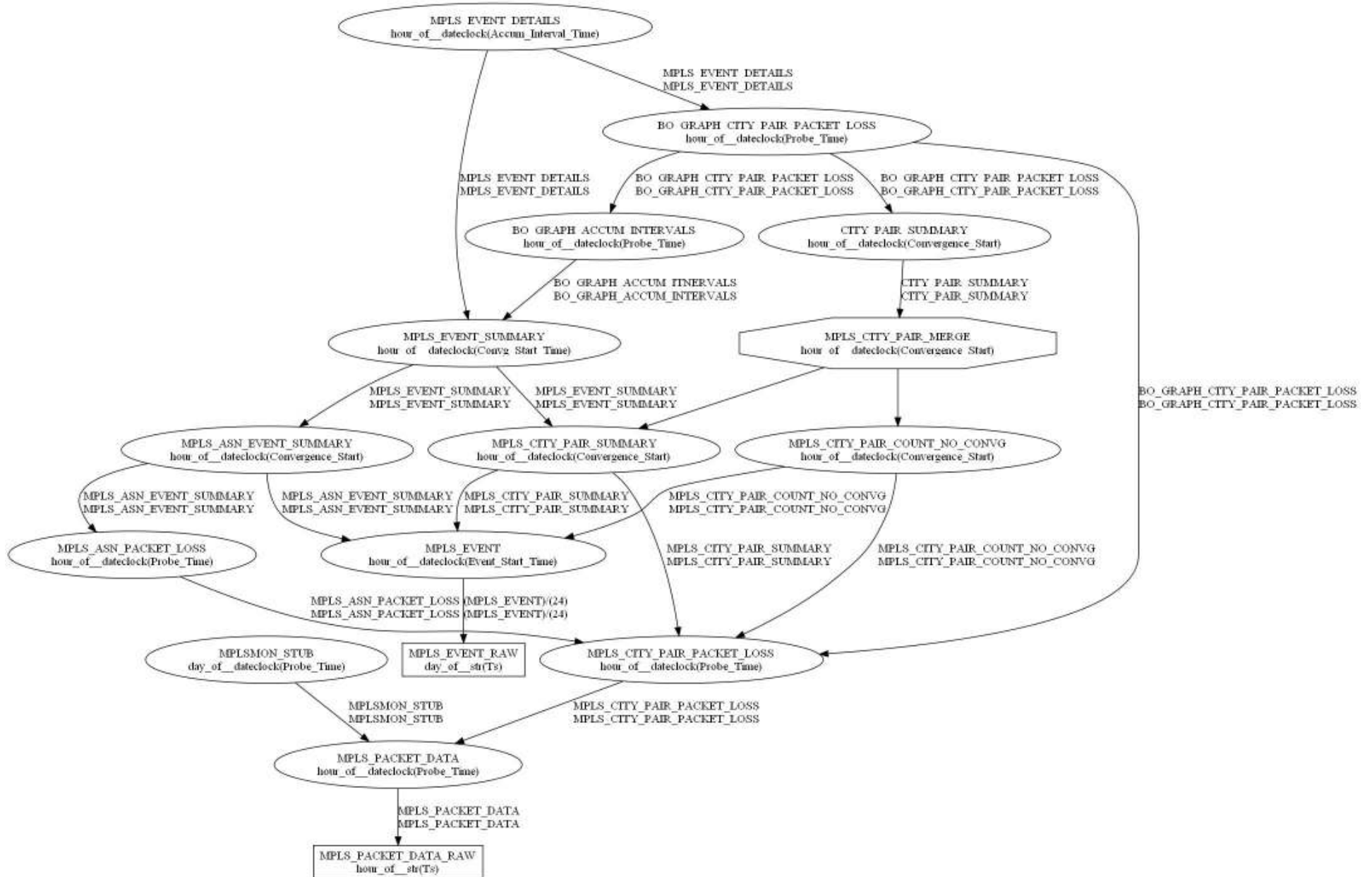
Darkstar: Backbone Network Data.



BPS



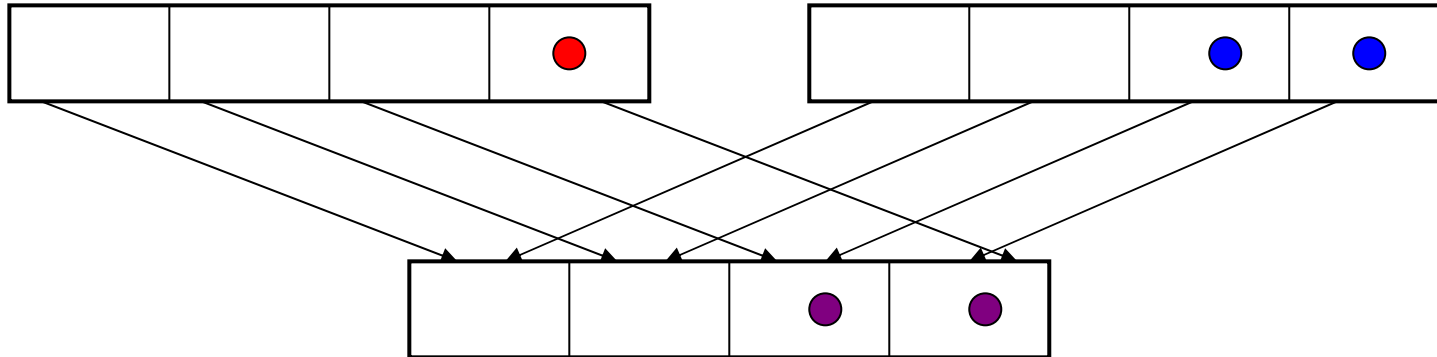
MPLS_MON



Update Propagation

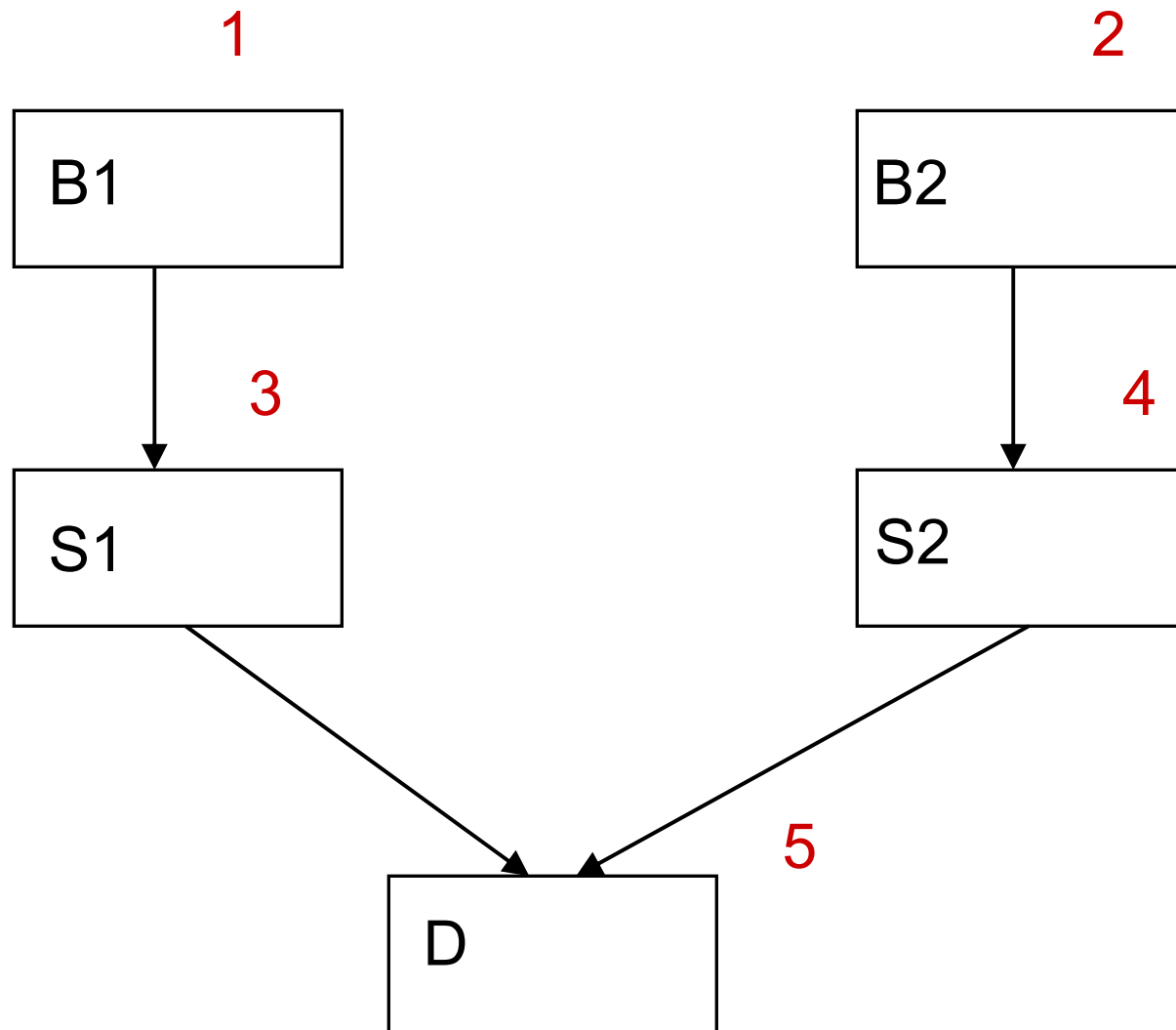
- We can build complex apps if we're confident that *all* updates get propagated.
- 1st version: used *make-style* algorithm
 - Not correct for complex configurations
 - Requires global analysis and all-at-once updates.
- Developed update propagation theory
 - Interaction with scheduler
 - Merge tables, partition rollups, etc.
- Implemented correct algorithm
 - Has some scheduling restrictions
- Problem: scheduling restrictions cause update delays in RT tables.
 - Move to algorithm without scheduler restrictions.

Incremental Updates

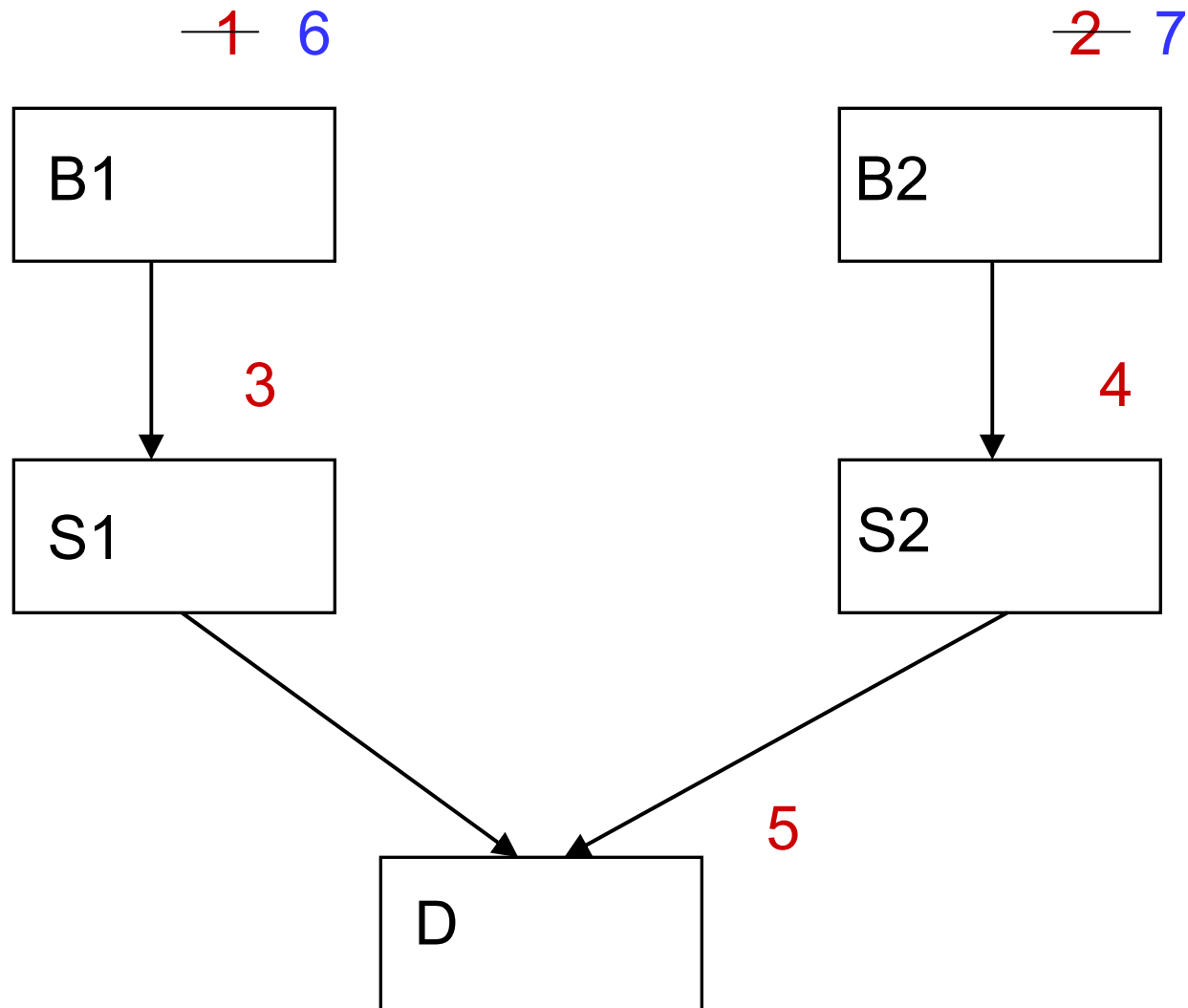


- Only propagate the increment.
- Update only those partitions whose sources have new data.
- How can we determine if a source partition has more recent data?

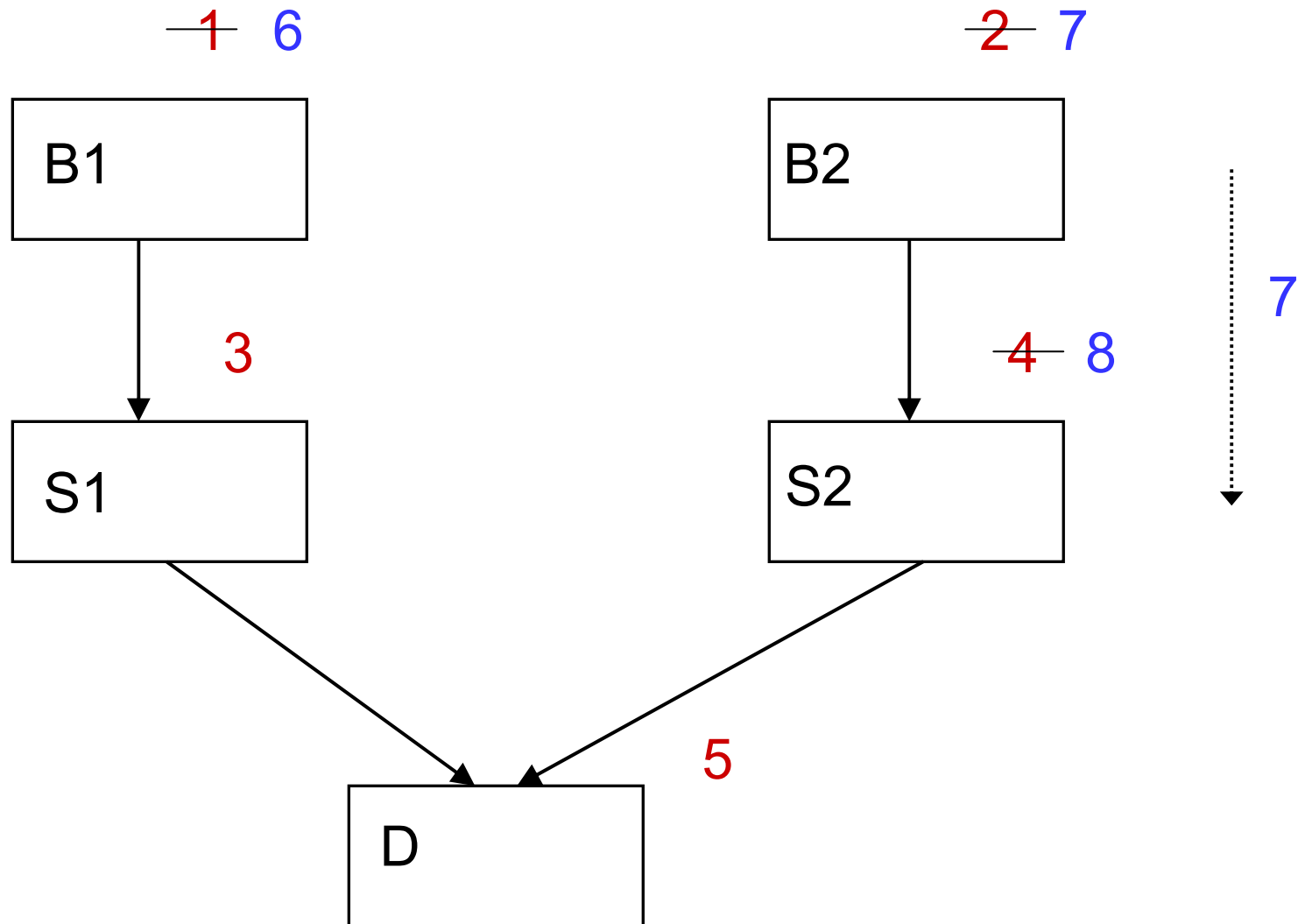
“make” doesn't work



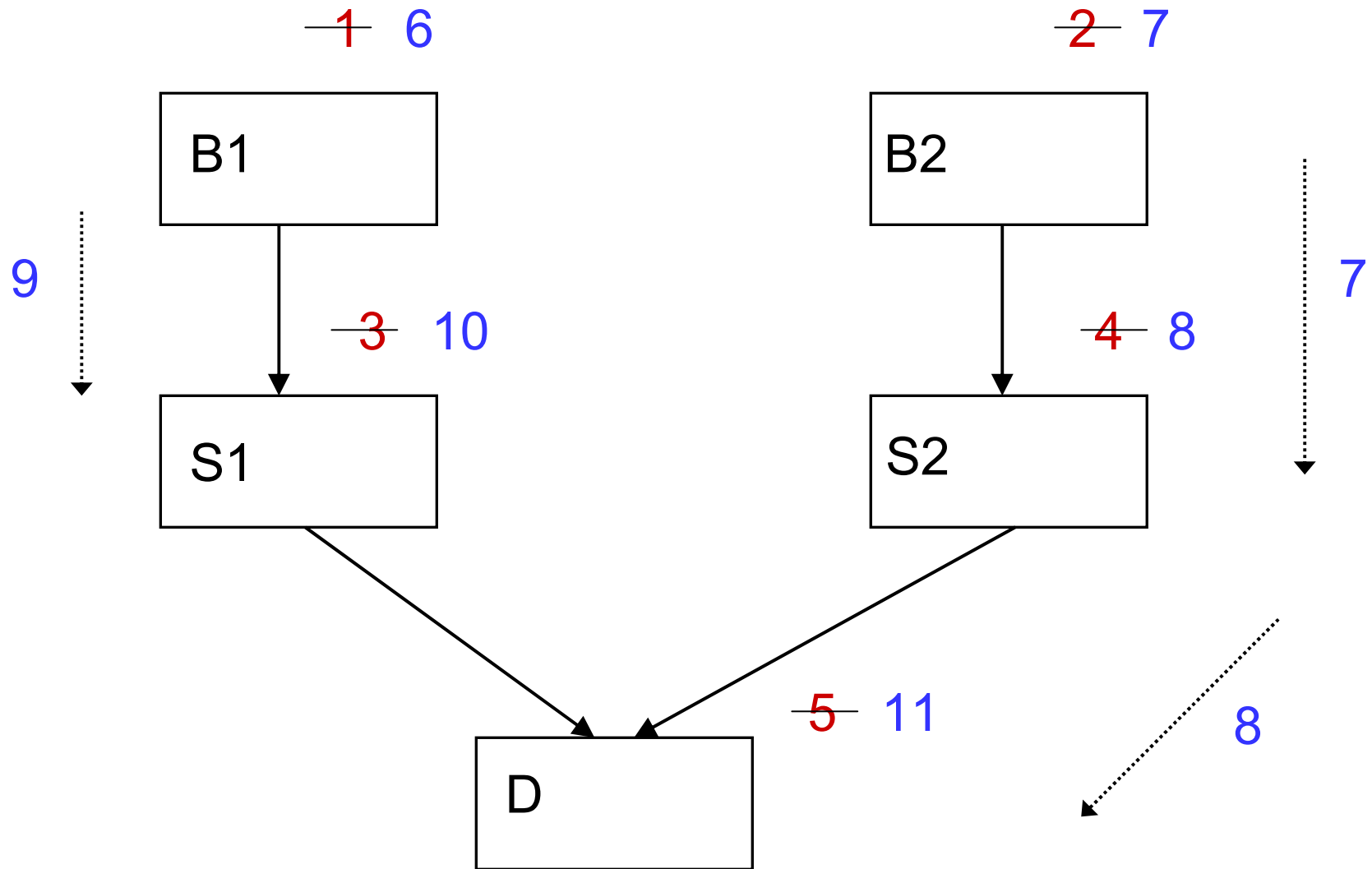
“make” doesn’t work



“make” doesn’t work

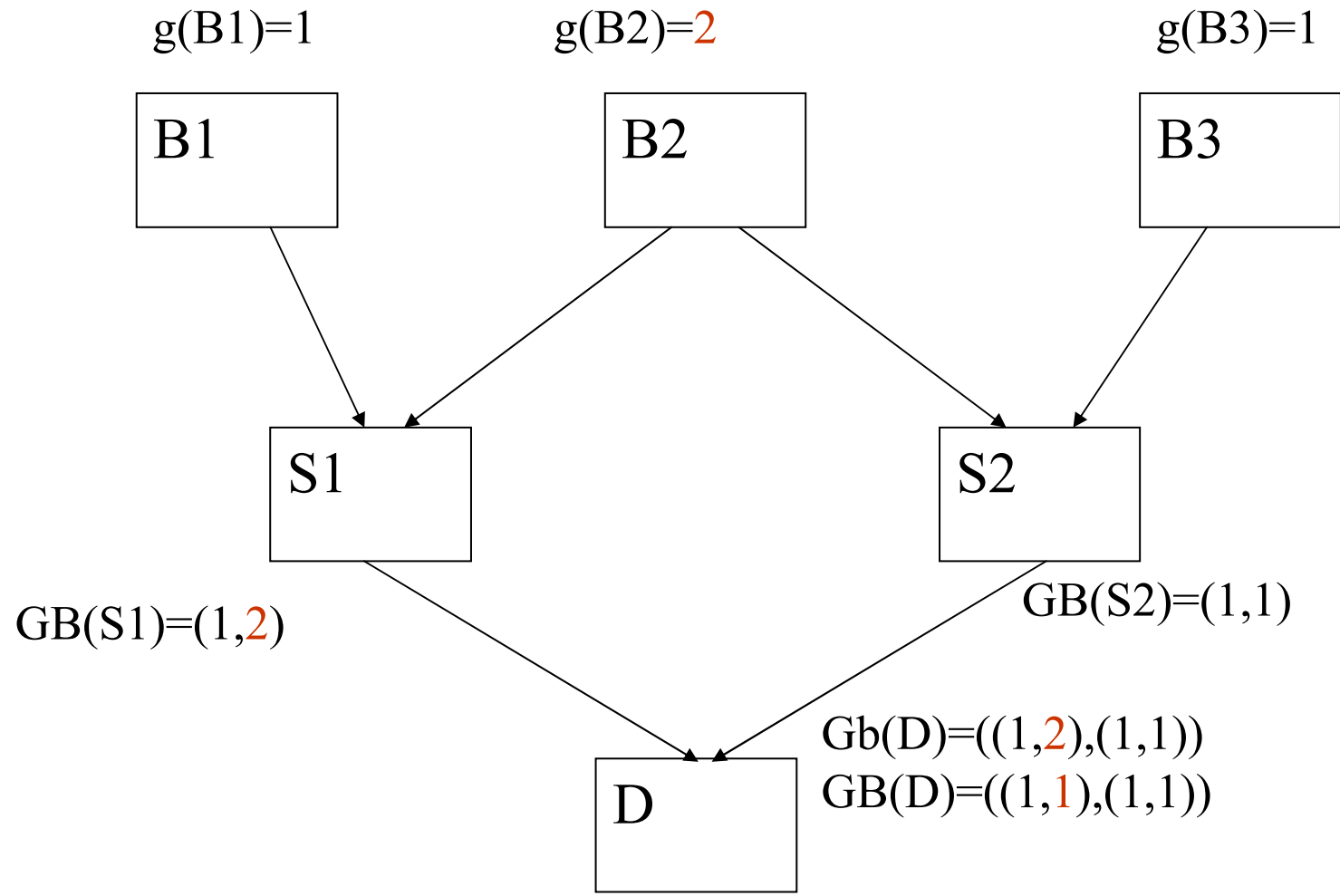


“make” doesn’t work



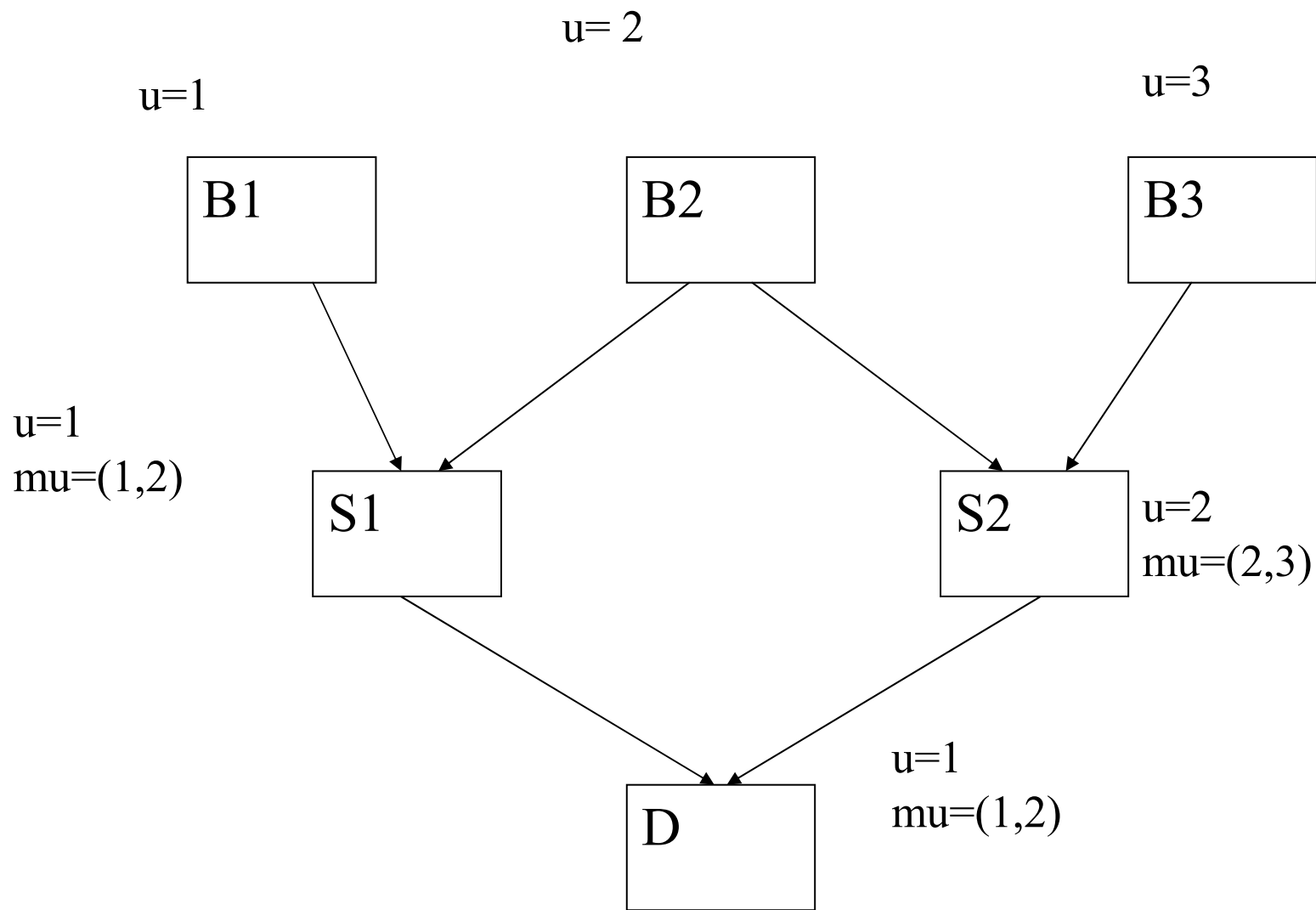
Correctness Theory

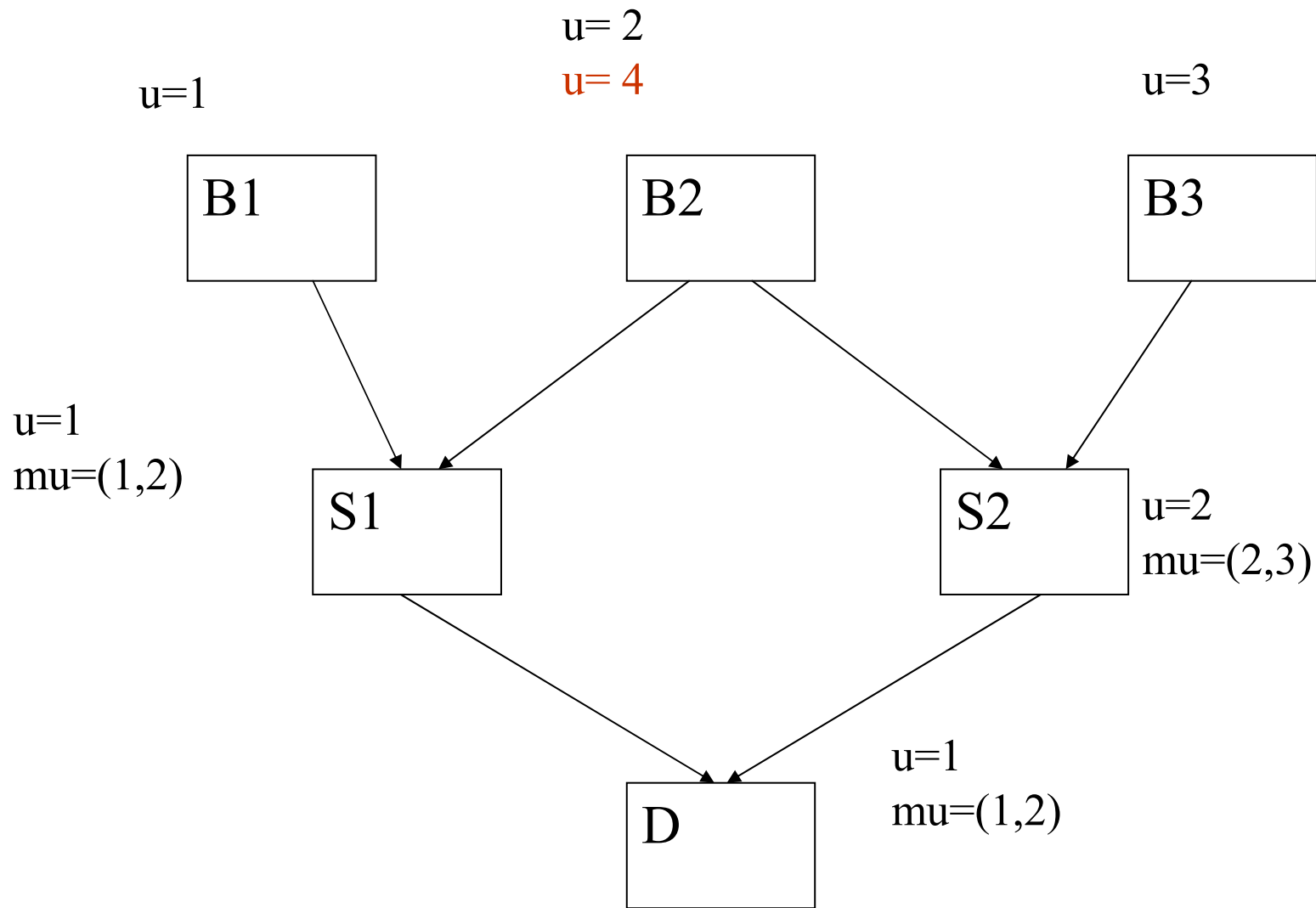
- The warehouse consists of *base* tables and *derived* tables.
- Each partition of a base table has a *generation number*: the number of times it has been updated.
- A derived table partition **D** depends on a collection **B(D)** of base table partitions
 - One entry for each distinct dependence path.
- **GB(D)** is the collection of generation numbers of the base table partitions in **B(D)** used to compute **D**.
- **Gb(D)** is the collection of base table partition numbers of **D**'s source partitions
- Update **D** iff. **Gb(D) > GB(D)**.

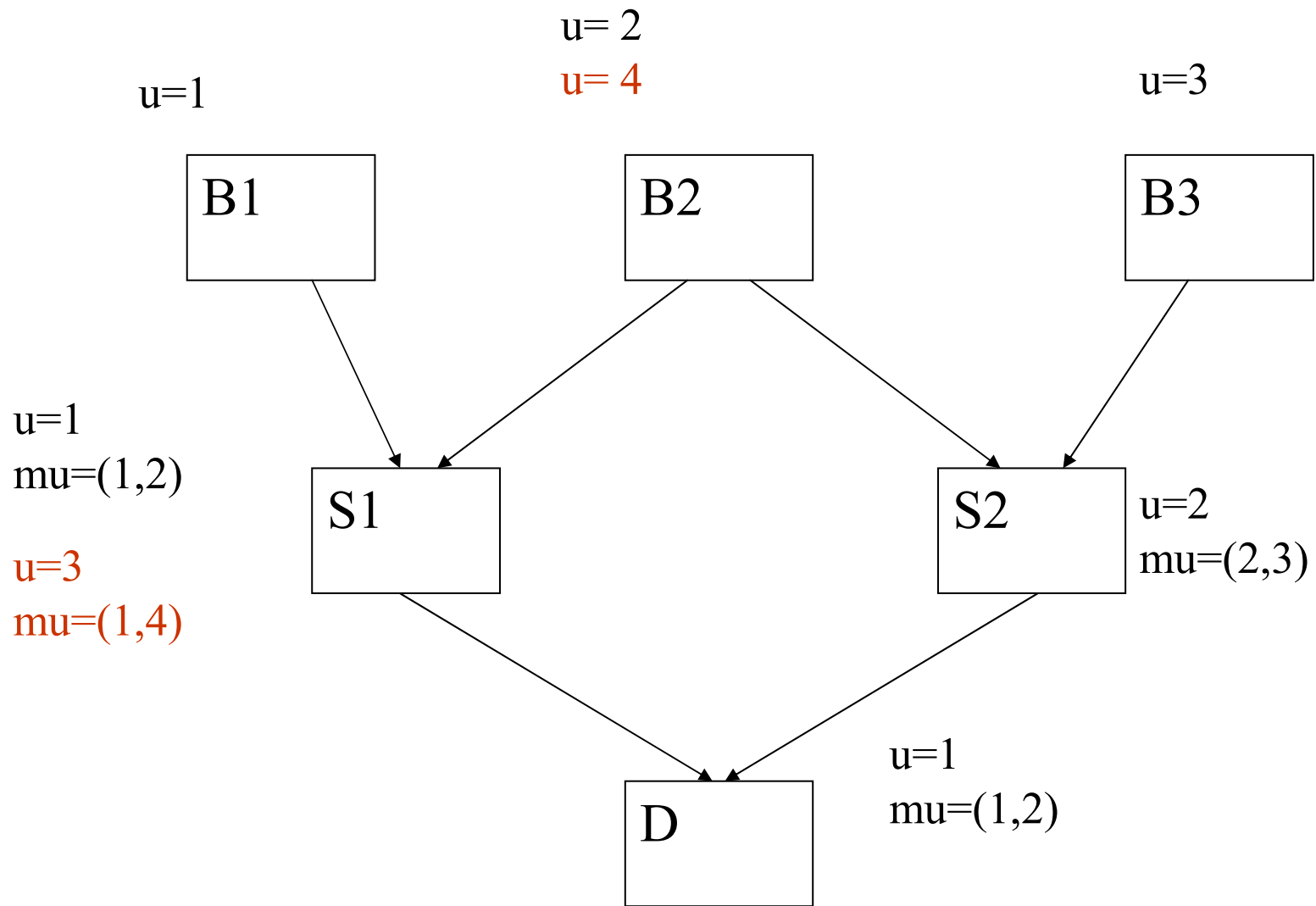


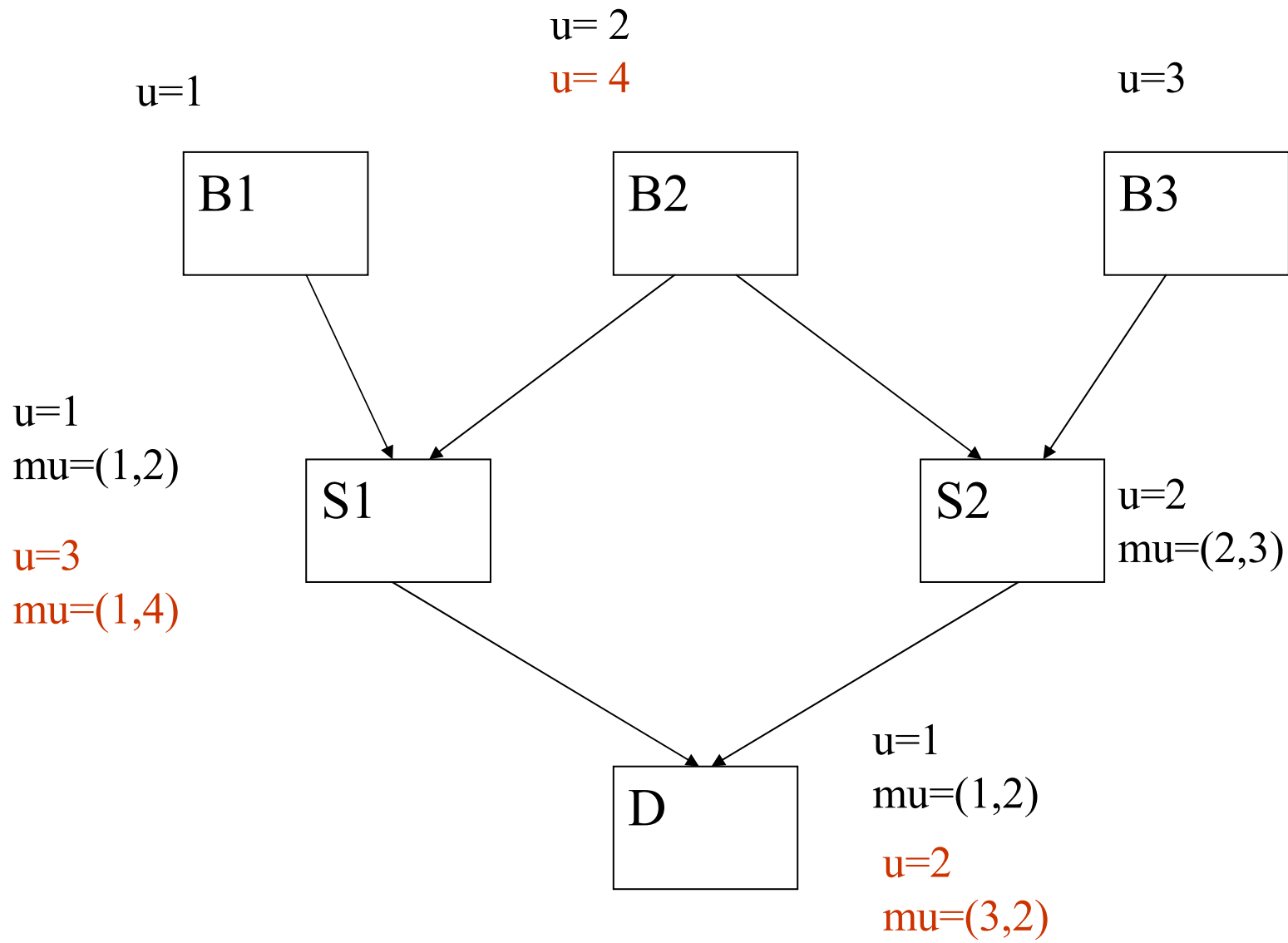
Source-vector Protocol

- Each table D maintains an update counter ut .
- Each partition d of D maintains an updatestamp u which is, assigned the value of ut when the partition is updated.
- Each partition also maintains a *source vector* $mu=(M1,...,Mn)$, which is the maximum updatestamp in the partitions from table S_i that affect partition d .
- Update partition d iff. there is a source table S_i with partitions $(s_{i1},...,s_{ij})$ such that $\max(u(s_{i1}),...,u(s_{ij})) > M_i$.

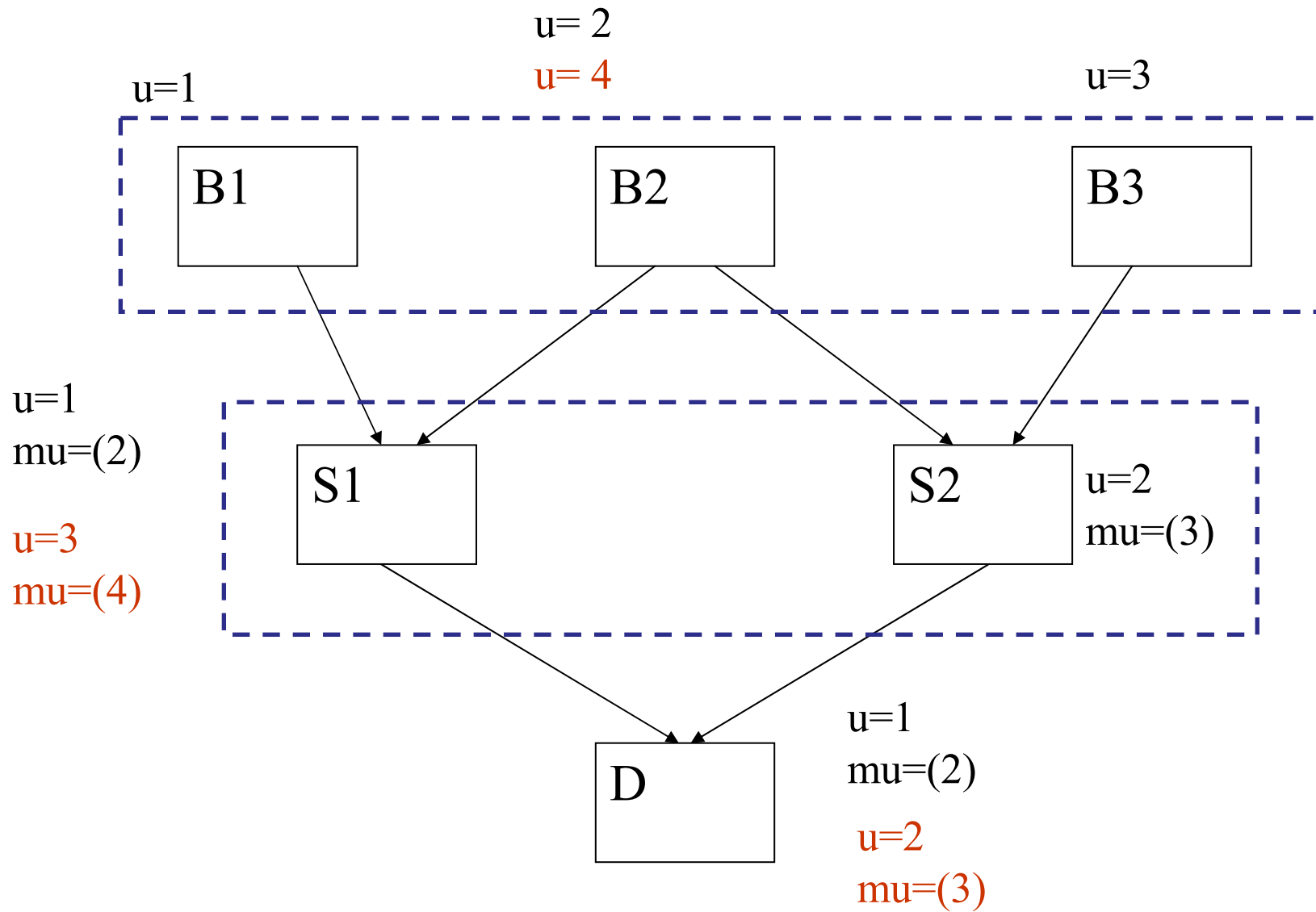








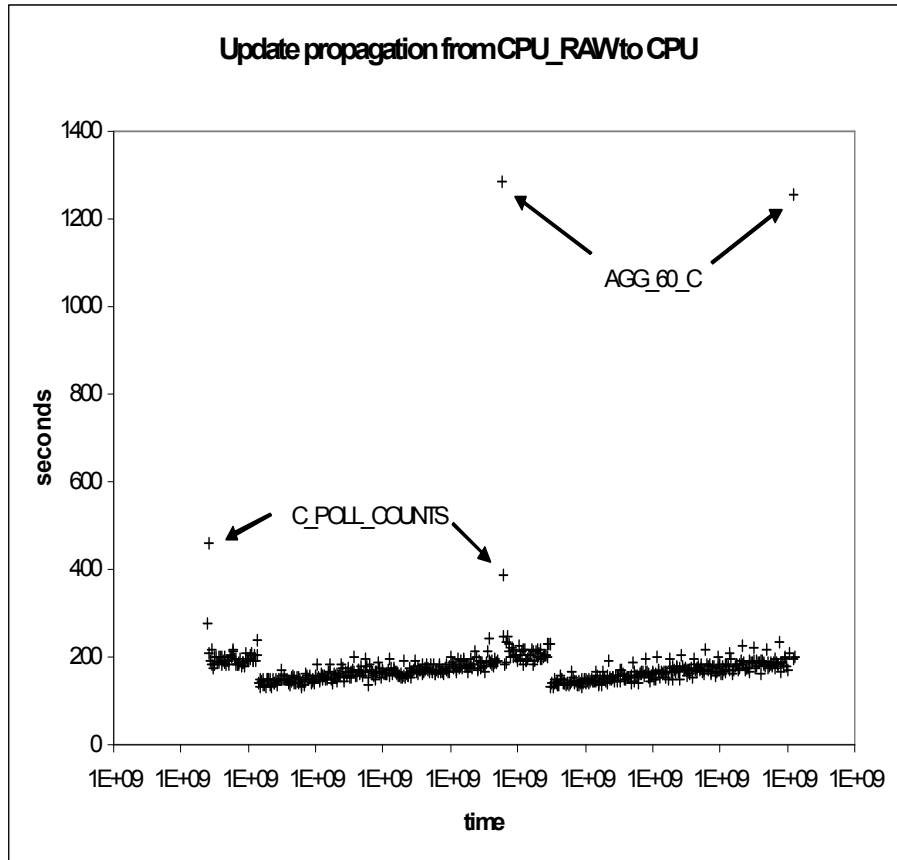
Suppose B1, B2, B3 are in table B, S1, S2 in table S



Extensions

- **Interval-timestamp protocol**
 - Store start and end timestamp of the update.
 - Fixed-size metadata, but it requires synchronized timestamps in a cluster
- **Trailing-edge consistency**
 - Propagate “no more updates” punctuation
- **MERGE tables (Union)**
 - Treat new partitions as newly loaded.
- **Partition Rollups**
 - Problem: load data once per minute, store for 2 years
 - 1-minute partitions => 1,051,200 partitions to cover a year
 - Rollup older partitions
 - First 2 days : 1 minute partitions
 - Next 728 days : 1 day partitions
 - 3608 partitions to cover 2 years, but with efficient real-time updates.
- **Any single-timestamp protocol requires scheduling restrictions**
 - Scheduling restrictions are bad for a real-time warehouse.

Effects of Scheduling Restrictions



Conclusions

- Update propagation is a critical component of an on-line stream warehouse.
 - Continual loading
 - Thousands of tables.
- Make algorithm doesn't work.
- Best algorithm uses a modified vector timestamp.
 - Aggressive timestamp compression.
- Critical component of many stream warehouses in use in AT&T.