# A Prolegomenon for Predictive Modeling of Parallel OLTP Systems

Andrew Pavlo Saurya Velagapudi Stanley Zdonik **Brown University** 

Evan P.C. Jones Samuel Madden **MIT CSAIL** 

# **Project Goal:**

For an OLTP workload containing distributed transactions, one would like to predict when transactions are not singlepartitioned. We are exploring optimizations for these types of workloads for the H-Store parallel, main-memory DBMS using partitioned predictive Markov models.

Client

# **OLTP Markov Models:**

Using the workload trace, the modeler generates a separate Markov model for each stored procedure. Each graph is specific to the initial partition in the database that can invoke the stored procedure. The vertices represent the potential execution states of a transaction, including (1) the executed query and (2) what partitions have been read/written up to that point. The edges weights represent the probabilities of one state transitioning to another.

#### **Single-Partition Transaction Multi-Partition Transaction**



No locks, No 2PC, No network

### Inputs:

#### Stored Procedures

H-Store requires that administrators predefine stored procedures comprised of Java control code intermixed with parameterized SQL statements.

#### Sample Workload Trace

The administrator must also provide a sample workload trace for the target application. A trace record consists of (1) the stored procedure executed (2) the queries executed in the transaction, and

 $\otimes$  locks,  $\otimes$  2PC,  $\otimes$  network callOrderLine **extends** StoredProcedu final String GetOrders = FROM ORDERS" + WHERE ((0\_ALL\_LOCAL = 1 AND 0\_W\_ID = ?)" +
OR (0\_ALL\_LOCAL = 0)) AND 0\_ENTRY\_D >= ?"; Lic final String GetOrderLines = ERE OL\_O\_ID = ? AND OL\_W\_ID = ?" + SUPPLY\_W\_ID = ? AND OL\_I\_ID = ?" -DELIVERY\_ID IS NULL": final String UpdateCustomer ER SET C\_BALANCE -= ?" HERE C\_ID = ? AND C\_W\_ID = ?"; Lic final String UpdateStock = PDATE STOCK SET S\_QUANTITY = 0 " + NHERE S\_W\_ID = ? AND S\_I ID = ?": ublic void run(long s\_w\_id, long i\_id, Date o\_date) Table orders = executeSQL(GetOrders, s\_w\_id, o\_date) for (Row record : orders.getRows()) long o\_id = record.get("0 ID"); long o\_w\_id = record.get("0\_W\_ID") Long o\_c\_id = record.get("0\_C\_ID"); Table order\_lines = executeSQL(GetOrderLines, o\_id, o\_w\_id, s\_w\_id, i\_id) **louble** refund = 0.0d; for (Row ol : order\_lines.getRows()) {
 refund += ol.get("OL\_QTY") \* ol.get("OL\_AMT"); executeSQL(UpdateCustomer, refund, o\_c\_id, o\_w\_id); executeSQL(UpdateStock, s\_w\_id, i\_id); RecallOrderLine (1001,123,2010-01-18) exec GetOrders (1001,2010-01-18) exec GetOrderLines (555,1002,1001,123) UpdateCustomer(\$10,19811,1002) UpdateStock (1001, 123) RecallOrderLine (1001,123,2010-01-18) exec GetOrders (1001,2010-01-18) exec GetOrderLines (555,1002,1001,123) exec UpdateCustomer (\$10,19811,1002) exec GetOrderLines (555,1002,1001,123) exec UpdateCustomer (\$10,10011,1002)

UpdateCustomer(\$10,19811,1002) UpdateStock (1001, 123)



#### (3) the input parameters.



# **Runtime Optimizations:**

An OLTP client sends a transaction request to a single H-Store execution node in the cluster.



The node's transaction coordinator invokes its procedure executor and estimates the execution path of the new transaction.



The estimator tracks state and sends the coordinator the relevant partition ids. The coordinator informs others on whether it is safe to execute new transactions on those partitions before the original transaction commits.



When the transaction commits or aborts, the node recomputes edge weights if the transaction deviated from the original predictions.



