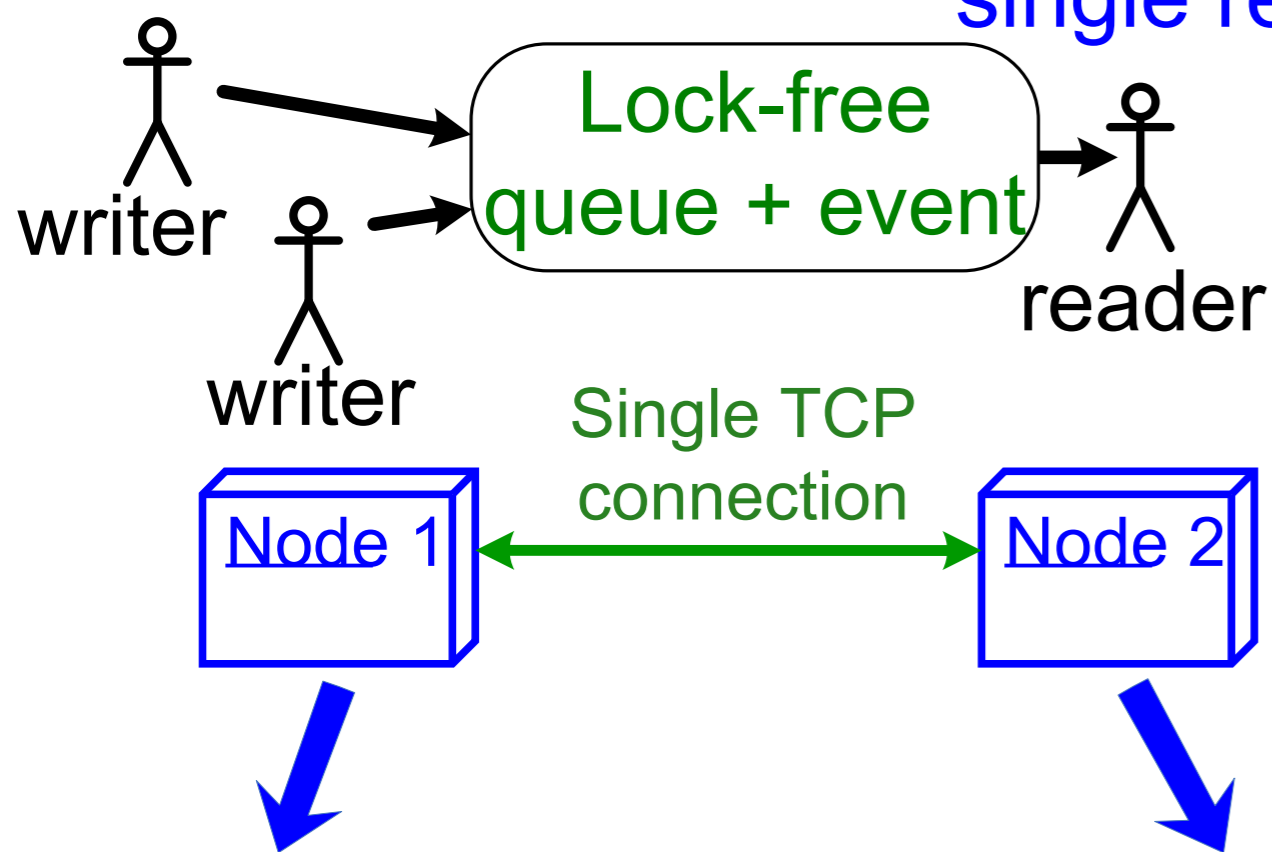




Saint-Petersburg University, Russia

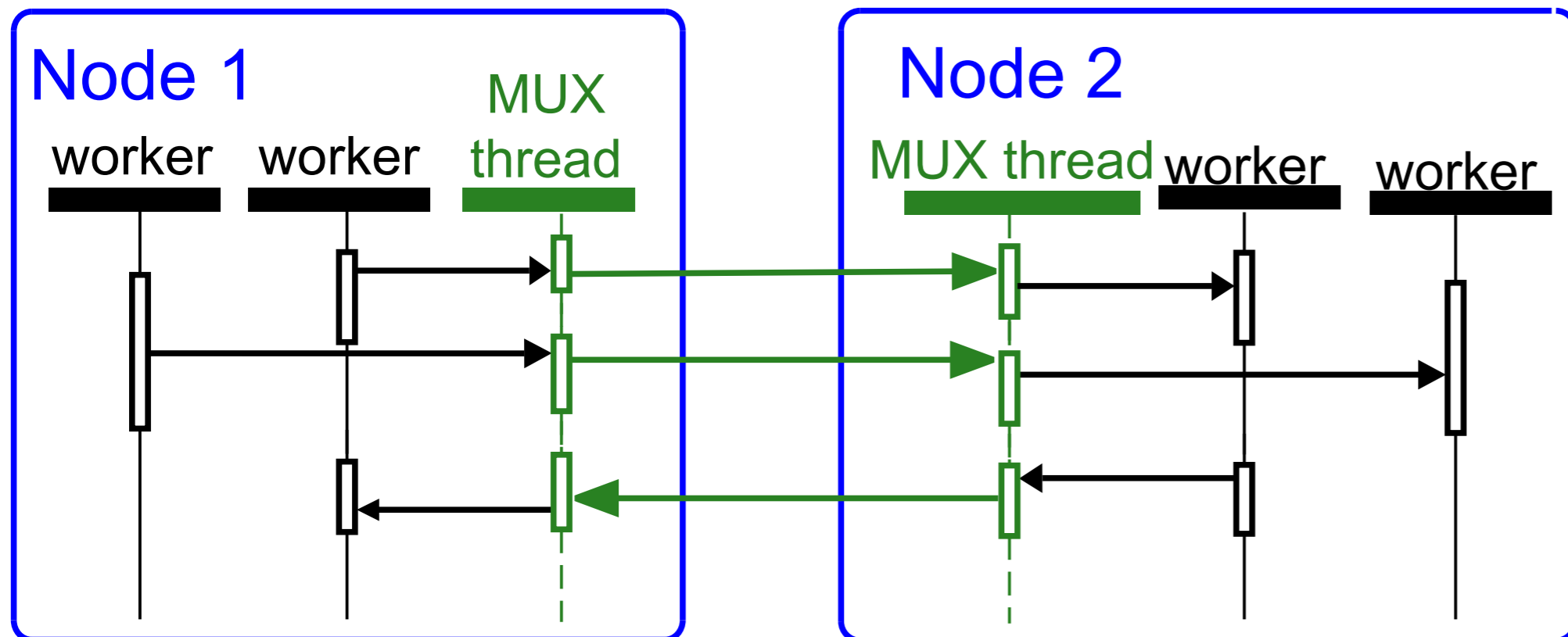
Team SPBU: George Chernishev, Kirill Smirnov

Multiple writers (50)



Event-driven single reader

- Single TCP connection
- Thread pool for indexed lookups
- Lock-free interaction
- Socket/event MUX



High-Level Aspects

-
- Pipelined query execution
 - Master node is responsible for complex operations
 - Slave nodes perform only scans and projections
- #### Local Reduction
- Predicate clean-up
 - Restriction propagation

Join Operation

- In-operation caching
- Sort one result, probe with another
- Tend to sort full result (to sort once)

Plan Generation

- Try to push down selects with indexed field
- Join indexed field with indexed (to fit in op cache)

Heuristics

- predicate pushdown
- operator selection

Statistics

- value distribution
- min/max
- duplicates

Partitioning

- number of partitions
- involved nodes

SELECT .. FROM .. WHERE

Search for "good" query plans

several candidates

Cost-based plan selection

best plan

Plan distribution and execution

Master Node

- translate
- recv
- query processing
- translation

Slave Node

- send
- join
- access
- join
- access
- recv
- partial plan execution
- data access
- join processing

Slave Node

- send
- access

Translator

- +next(Tuple)
- +open()
- +close()
- +reset()

Vectorized volcano-style physical operators

Translator adapts vectorized to one-tuple-a-time model

Operator

- +next(Page)
- +open()
- +close()
- +reset()

- Block NestedLoop Join
default if no other join fits
- PartitionedHash Join
uses the hash-based index
- Merge Join
if already sorted

TCP sockets tuned for performance: NODELAY, TCP_CORK

- TBScan
plain old table scanner
- TIDGrabber
fetches rows from file by Tuple Identifier
- IXScan
evaluates a single index
- MultiIXScan
handles multiple indexes in a single run

Black Arts

a lot of nasty compile-time optimizations

Future steps

- move from read-only to write-supported DBMS
- recycle code-base for future project: an Energy-Proportional DBMS running on multiple nodes
- dynamic node-switching wrt. performance needs
- turn off idle nodes for power savings
- build a scalable, energy-aware database cluster

Daniel Schall
Volker Hudlet

Database and Information Systems
University of Kaiserslautern

{schall, hudlet}@cs.uni-kl.de
http://www.lgis.cs.uni-kl.de

Contact

Supported by

