

Becoming friends

— A common foundation for data processing —

Sebastian Ertel

Programming Languages vs. Databases

Analytical Queries
(TPC-H)

Data Processing Engine

Transactional Programs
(TPC-C)

Data Integration
(ETL, BI-Tools)

Data Streaming
(CQL, etc.)



Data Mining
(Machine Learning)

Programming Languages vs. Databases

Analytical Queries
(TPC-H)

Data Processing Engine

Transactional Programs
(TPC-C)

Data Integration
(ETL, BI-Tools)

Data Streaming
(CQL, etc.)



Data Mining
(Machine Learning)

Programming Languages vs. Databases

Analytical Queries
(TPC-H)

SQL / JavaScript

Data Processing Engine

C/C++, Java, (Scala)

Transactional Programs
(TPC-C)

SQL / C/C++, Java, etc.

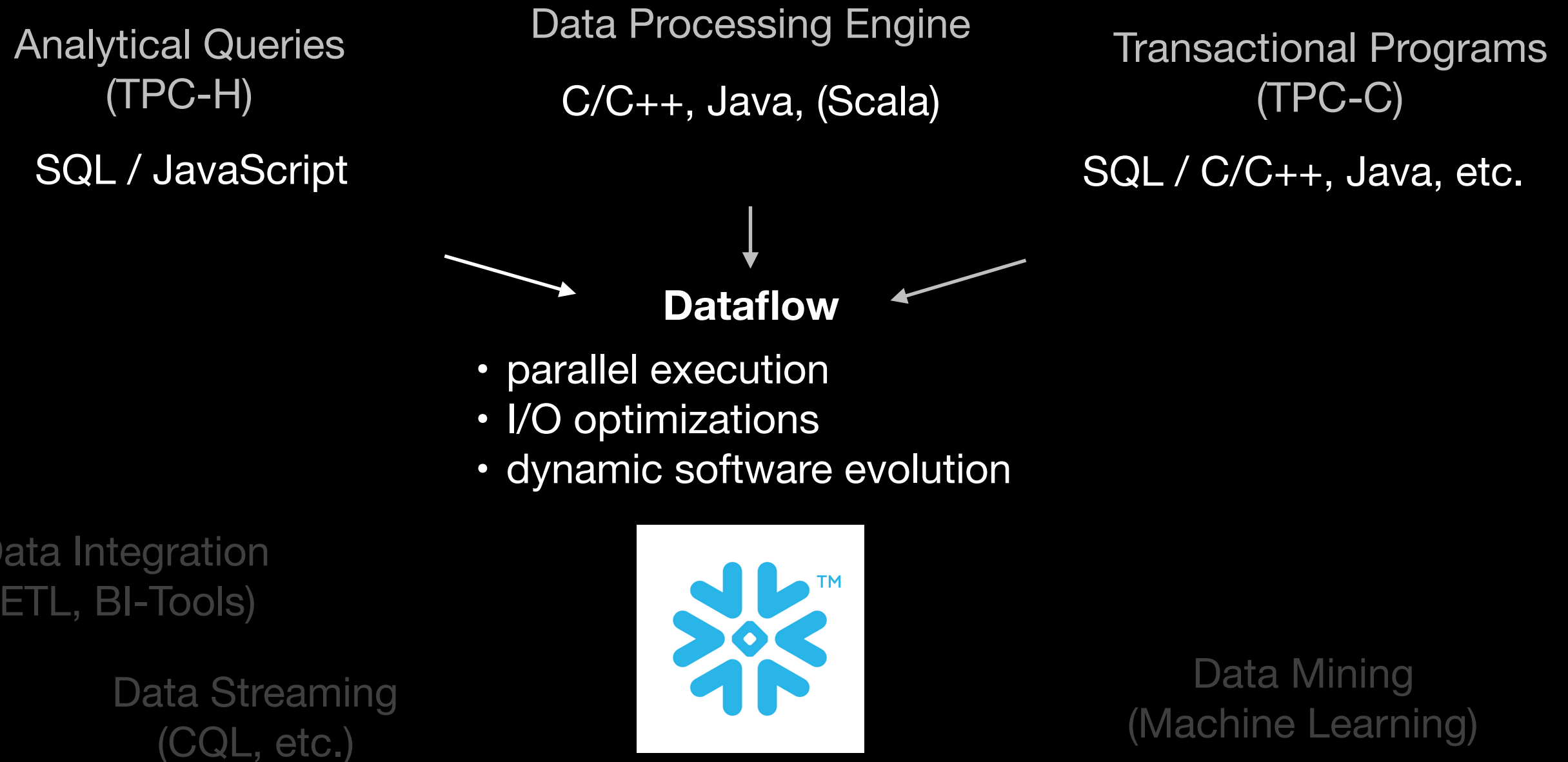
Data Integration
(ETL, BI-Tools)

Data Streaming
(CQL, etc.)



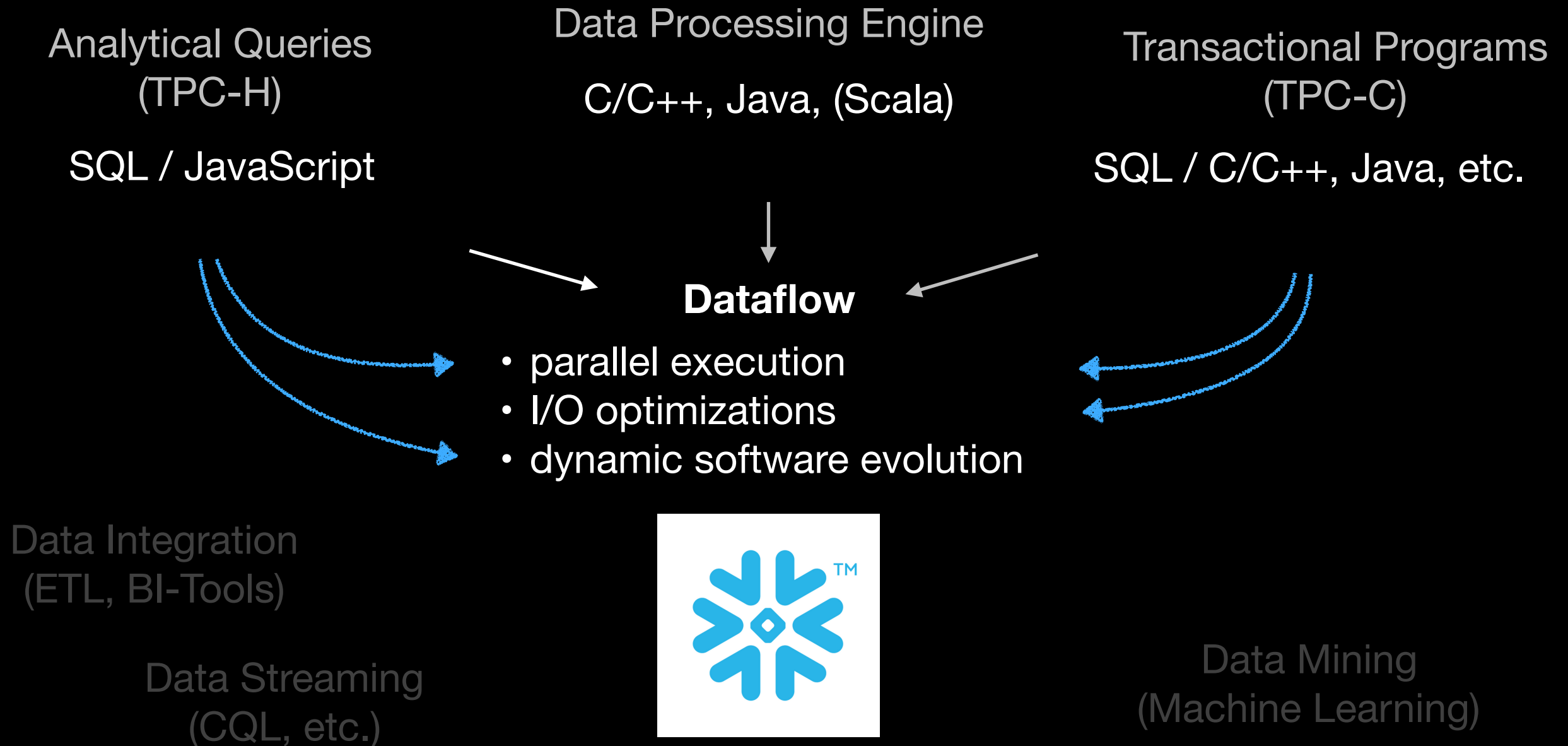
Data Mining
(Machine Learning)

Programming Languages vs. Databases



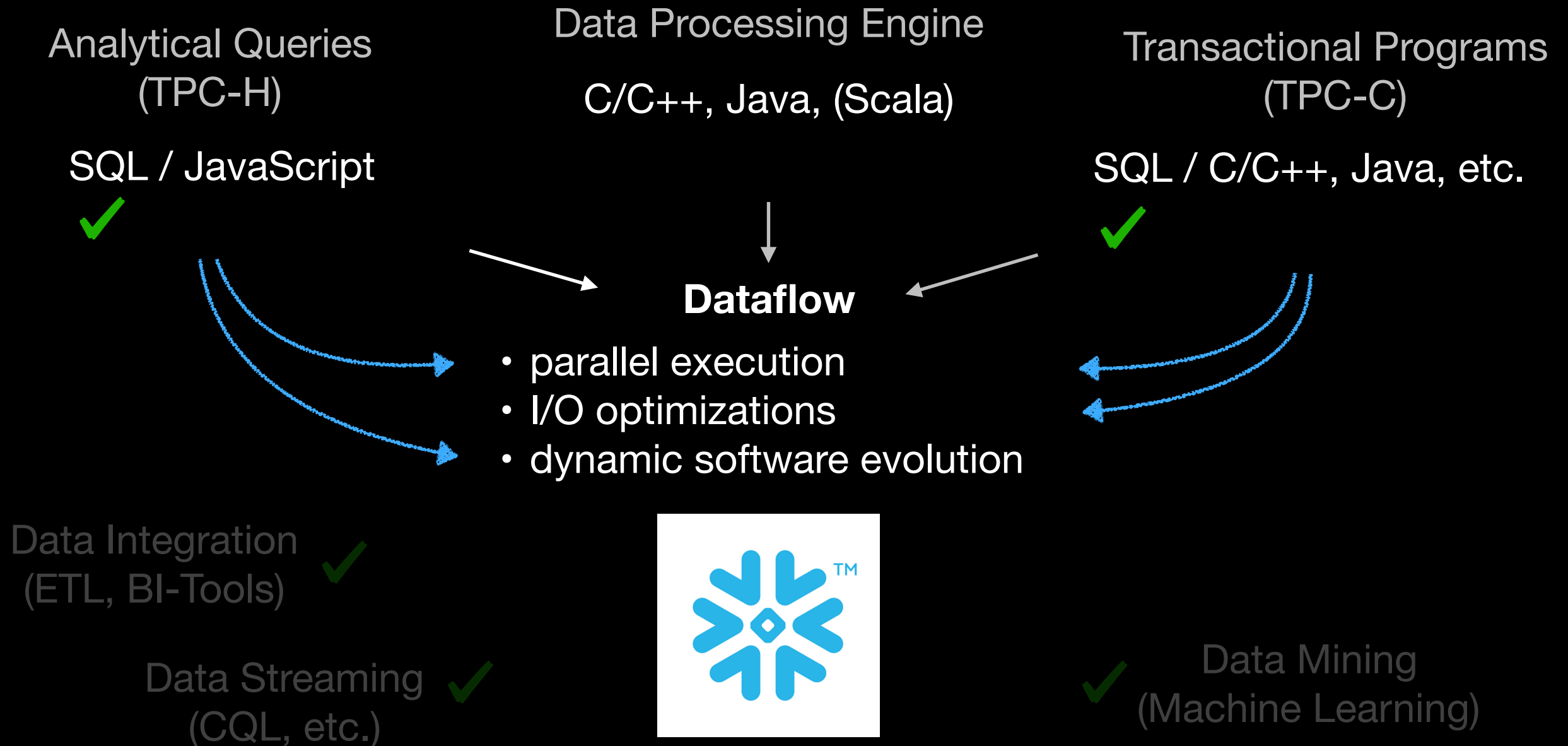
Sebastian Ertel, Andrés Goens, Justus Adam, and Jeronimo Castrillon. *Compiling for concise code and efficient I/O*. CC 2018.
Sebastian Ertel, Justus Adam, and Jeronimo Castrillon. *Supporting fine-grained dataflow parallelism in big data systems*. PMAM 2018.
Sebastian Ertel and Pascal Felber. *A framework for the dynamic evolution of highly-available data flow programs*. Middleware 2015.

Programming Languages vs. Databases



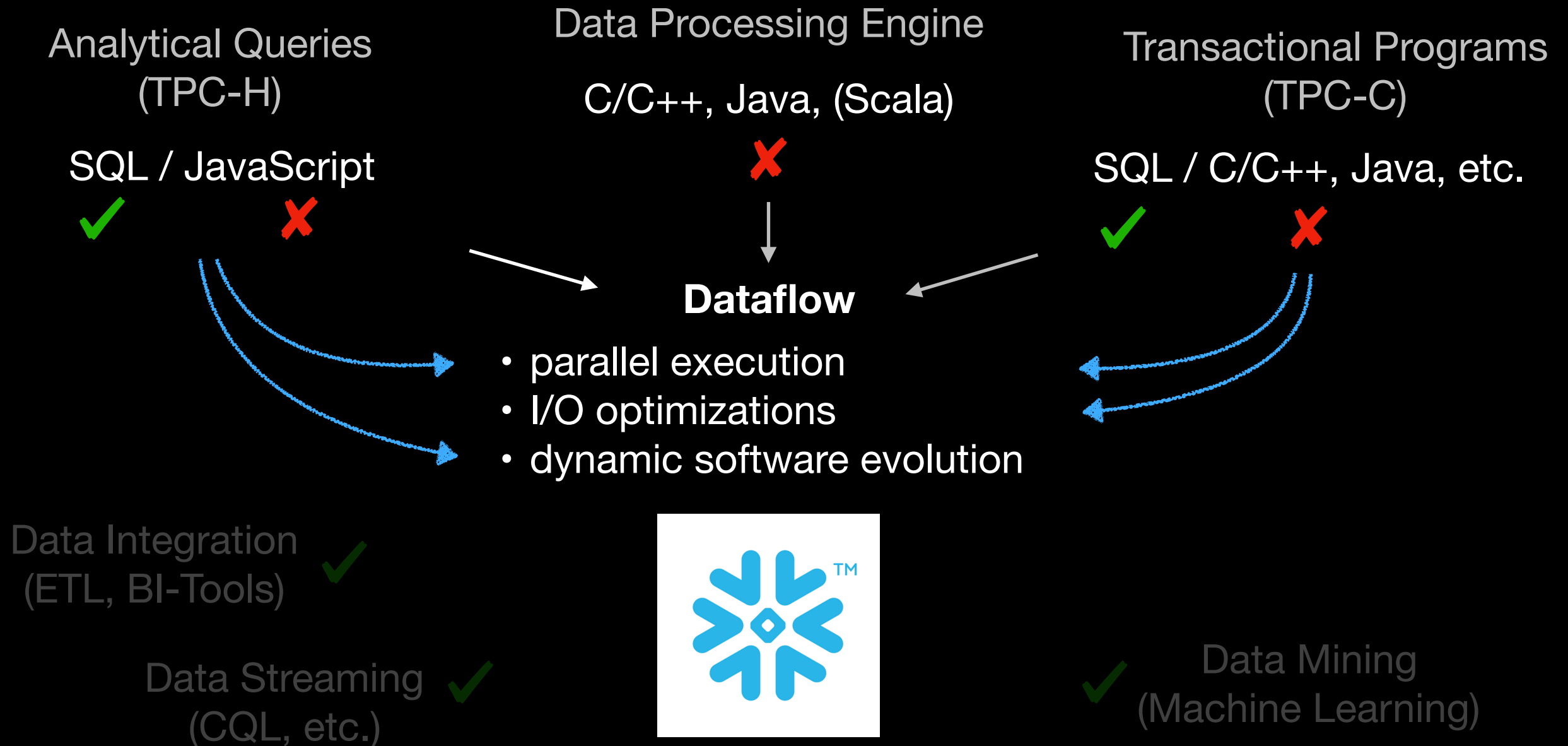
Sebastian Ertel, Andrés Goens, Justus Adam, and Jeronimo Castrillon. *Compiling for concise code and efficient I/O*. CC 2018.
Sebastian Ertel, Justus Adam, and Jeronimo Castrillon. *Supporting fine-grained dataflow parallelism in big data systems*. PMAM 2018.
Sebastian Ertel and Pascal Felber. *A framework for the dynamic evolution of highly-available dataflow programs*. Middleware 2015.

Programming Languages vs. Databases



Sebastian Ertel, Andrés Goens, Justus Adam, and Jeronimo Castrillon. *Compiling for concise code and efficient I/O*. CC 2018.
Sebastian Ertel, Justus Adam, and Jeronimo Castrillon. *Supporting fine-grained dataflow parallelism in big data systems*. PMAM 2018.
Sebastian Ertel and Pascal Felber. *A framework for the dynamic evolution of highly-available data flow programs*. Middleware 2015.

Programming Languages vs. Databases



Sebastian Ertel, Andrés Goens, Justus Adam, and Jeronimo Castrillon. *Compiling for concise code and efficient I/O*. CC 2018.
Sebastian Ertel, Justus Adam, and Jeronimo Castrillon. *Supporting fine-grained dataflow parallelism in big data systems*. PMAM 2018.
Sebastian Ertel and Pascal Felber. *A framework for the dynamic evolution of highly-available data flow programs*. Middleware 2015.

Enter hua!

Algorithm:

(Stateful) function:

Enter hua!

Algorithm:

```
void mapAlgo(InputIterator recordsOnDisk) {  
    var mapper = initMapper();  
    var writer = initWriter();  
  
    for(var chunk : recordsOnDisk){  
        var (line, content) = chunk;  
        var kvPairs = mapper.map(line, content);  
        for(var kvPair : kvPairs){  
            var (key, value) = kvPair;  
            writer.output(key, value);  
        }  
    }  
}
```

(Stateful) function:

Enter hua!

Algorithm:

```
void mapAlgo(InputIterator recordsOnDisk) {  
    var mapper = initMapper();  
    var writer = initWriter();  
  
    for(var chunk : recordsOnDisk){  
        var (line, content) = chunk;  
        var kvPairs = mapper.map(line, content);  
        for(var kvPair : kvPairs){  
            var (key, value) = kvPair;  
            writer.output(key, value);  
        }  
    }  
}
```

(Stateful) function:

```
Writer initWriter(){  
    return new Writer();  
}  
  
class Writer {  
    private Mapper$Context _ctxt; // state  
  
    public Writer(){ /* code omitted */ }  
  
    void write(Object key, Object value){  
        _ctxt.write(key, value);  
    }  
}
```

Enter hua!

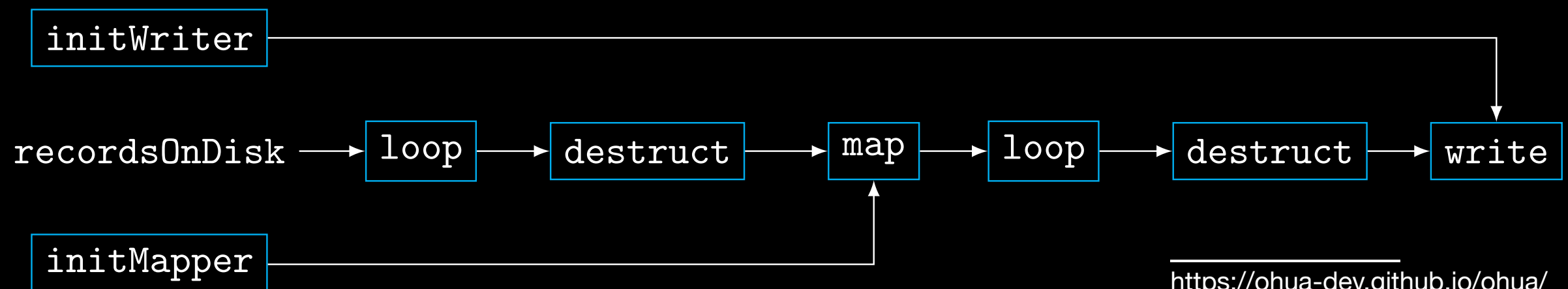
Algorithm:

```
void mapAlgo(InputIterator recordsOnDisk) {  
    var mapper = initMapper();  
    var writer = initWriter();  
  
    for(var chunk : recordsOnDisk){  
        var (line, content) = chunk;  
        var kvPairs = mapper.map(line, content);  
        for(var kvPair : kvPairs){  
            var (key, value) = kvPair;  
            writer.output(key, value);  
        }  
    }  
}
```

(Stateful) function:

```
Writer initWriter(){  
    return new Writer();  
}  
  
class Writer {  
    private Mapper$Context _ctxt; // state  
  
    public Writer(){ /* code omitted */ }  
  
    void write(Object key, Object value){  
        _ctxt.write(key, value);  
    }  
}
```

Dataflow Graph:



<https://ohua-dev.github.io/ohua/>

Enter Ohua!

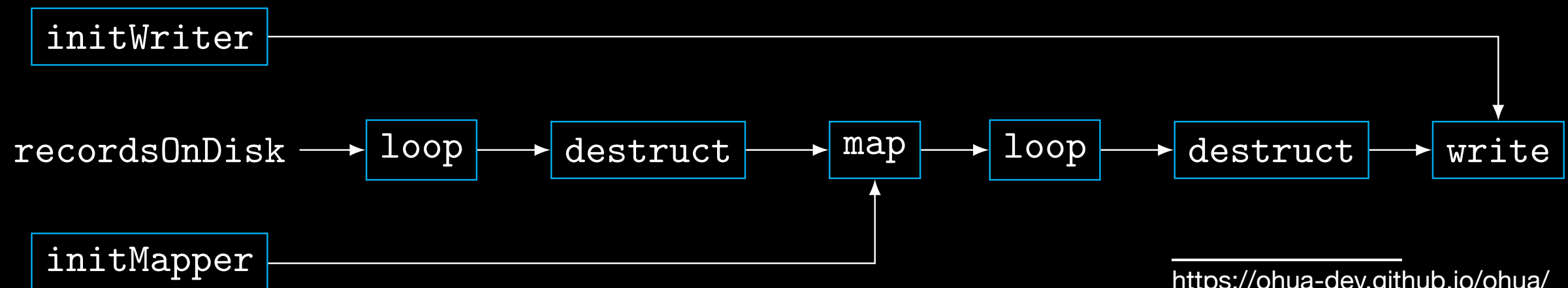
Algorithm:

```
void mapAlgo(InputIterator recordsOnDisk) {  
    var mapper = initMapper();  
    var writer = initWriter();  
  
    for(var chunk : recordsOnDisk){  
        var (line, content) = chunk;  
        var kvPairs = mapper.map(line, content);  
        for(var kvPair : kvPairs){  
            var (key, value) = kvPair;  
            writer.output(key, value);  
        }  
    }  
}
```

(Stateful) function:

```
Writer initWriter(){  
    return new Writer();  
}  
  
class lambda-calculus-based IR  
    private Mapper$Context _ctxt; // state  
  
    public Writer(){ /* code omitted */ }  
  
    void write(Object key, Object value){  
        _ctxt.write(key, value);  
    }  
}
```

Dataflow Graph:



Enter Ohua!

Algorithm:

```
void mapAlgo(InputIterator recordsOnDisk) {  
    var mapper = initMapper();  
    var writer = initWriter();  
  
    for(var chunk : recordsOnDisk) {  
        var (line, content) = chunk; // parse  
        var kvPairs = mapper.map(line, content);  
        for(var kvPair : kvPairs) {  
            var (key, value) = kvPair;  
            writer.output(key, value);  
        }  
    }  
}
```

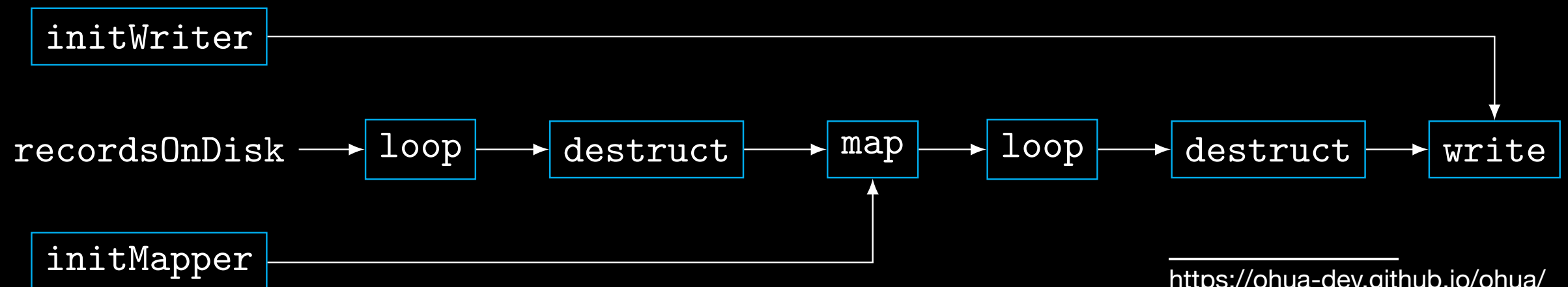
Rust
Java
Scala
C
C++
Haskell
JavaScript
...

(Stateful) function:

```
Writer initWriter(){  
    return new Writer();  
}  
  
class Mapper {  
    private Mapper$Context _ctxt; // state  
  
    public Writer() { /* code omitted */ }  
  
    void write(Object key, Object value){  
        _ctxt.write(key, value);  
    }  
}
```

lambda-calculus-based
IR

Dataflow Graph:



<https://ohua-dev.github.io/ohua/>

Enter Ohua!

Algorithm:

```
void mapAlgo(InputIterator recordsOnDisk) {  
    var mapper = initMapper();  
    var writer = initWriter();  
  
    for(var chunk : recordsOnDisk) {  
        var (line, content) = chunk;  
        var kvPairs = mapper.map(line, content);  
        for(var kvPair : kvPairs) {  
            var (key, value) = kvPair;  
            writer.output(key, value);  
        }  
    }  
}
```

Rust
Java
Scala
C
C++
Haskell
JavaScript
...



(Stateful) function:

lambda-calculus-based
IR

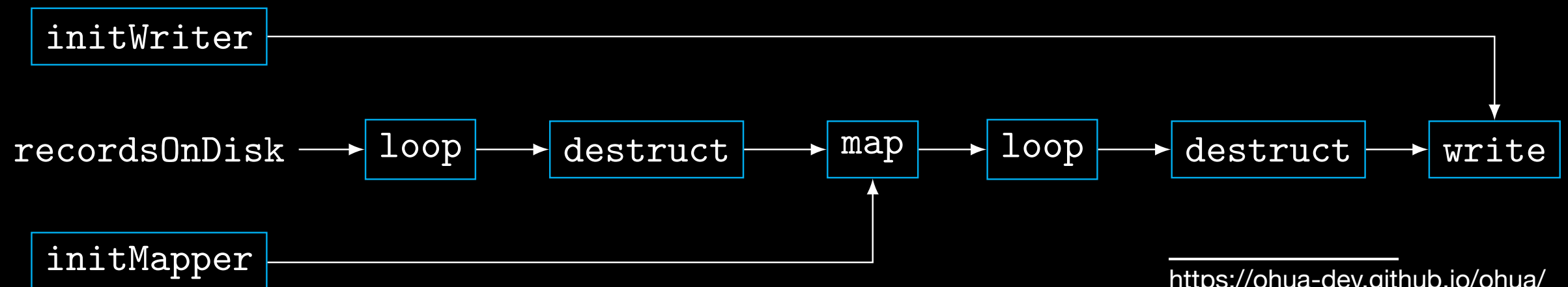
generic optimizations

domain-specific
optimizations

```
Write re  
re  
}  
clas  
pr  
; // state  
put  
tted */ }  
  
void write(Object key, Object value){  
    _ctxt.write(key, value);  
}  
}
```



Dataflow Graph:



<https://ohua-dev.github.io/ohua/>

Fully-integrated Dataflow Parallelism

Analytical Queries
(TPC-H)

SQL / JavaScript



Transactional Programs
(TPC-C)

SQL / C/C++, Java, etc.



Dataflow



Challenge: Speed is in SQL vs. rich JavaScript eco system!

Fully-integrated Dataflow Parallelism

Analytical Queries
(TPC-H)

SQL / JavaScript



Transactional Programs
(TPC-C)

SQL / C/C++, Java, etc.



Example UDF:

```
function factorial(D) {  
  if (D <= 0) {  
    return 1;  
  } else {  
    var result = 1;  
    for (var i = 2; i <= D; i++) {  
      result = result * i;  
    }  
    return result;  
  }  
}
```

Dataflow

Fully-integrated Dataflow Parallelism

Analytical Queries
(TPC-H)

SQL / JavaScript



Transactional Programs
(TPC-C)

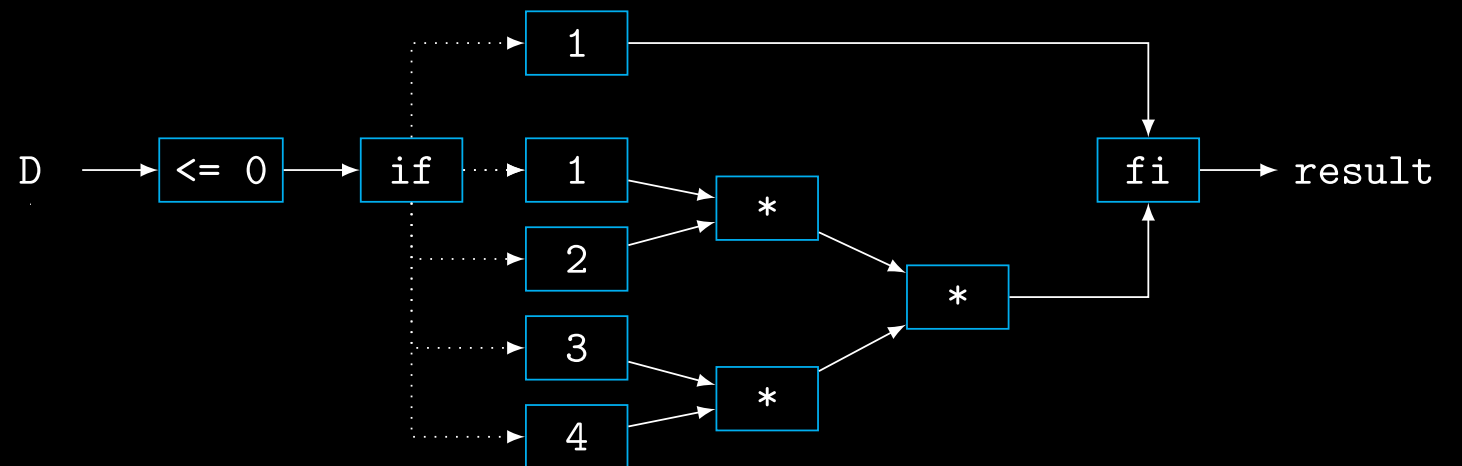
SQL / C/C++, Java, etc.



Example UDF:

```
function factorial(D) {  
  if (D <= 0) {  
    return 1;  
  } else {  
    var result = 1;  
    for (var i = 2; i <= D; i++) {  
      result = result * i;  
    }  
    return result;  
  }  
}
```

Dataflow



Fully-integrated Dataflow Parallelism

Analytical Queries
(TPC-H)

SQL / JavaScript



Transactional Programs
(TPC-C)

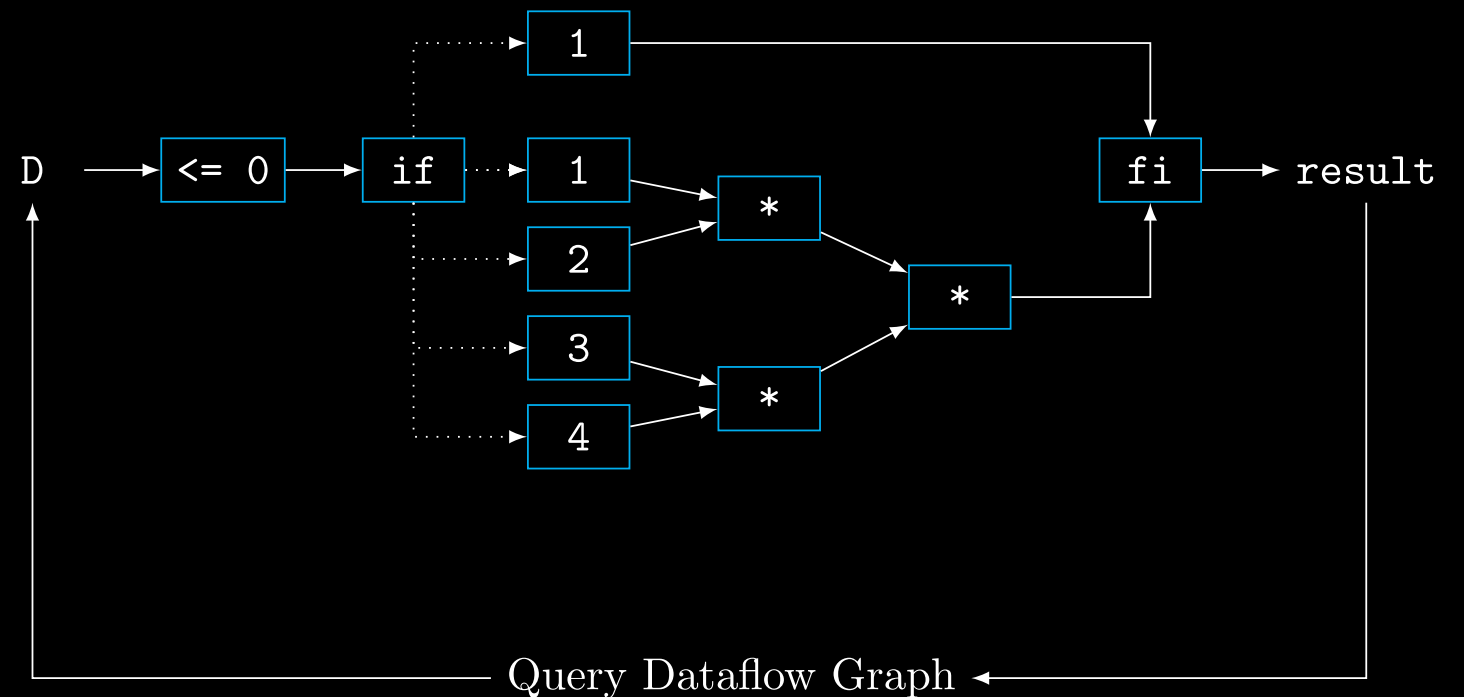
SQL / C/C++, Java, etc.



Example UDF:

```
function factorial(D) {  
  if (D <= 0) {  
    return 1;  
  } else {  
    var result = 1;  
    for (var i = 2; i <= D; i++) {  
      result = result * i;  
    }  
    return result;  
  }  
}
```

Dataflow



Fully-integrated Dataflow Parallelism

Analytical Queries
(TPC-H)

SQL / JavaScript



Transactional Programs
(TPC-C)

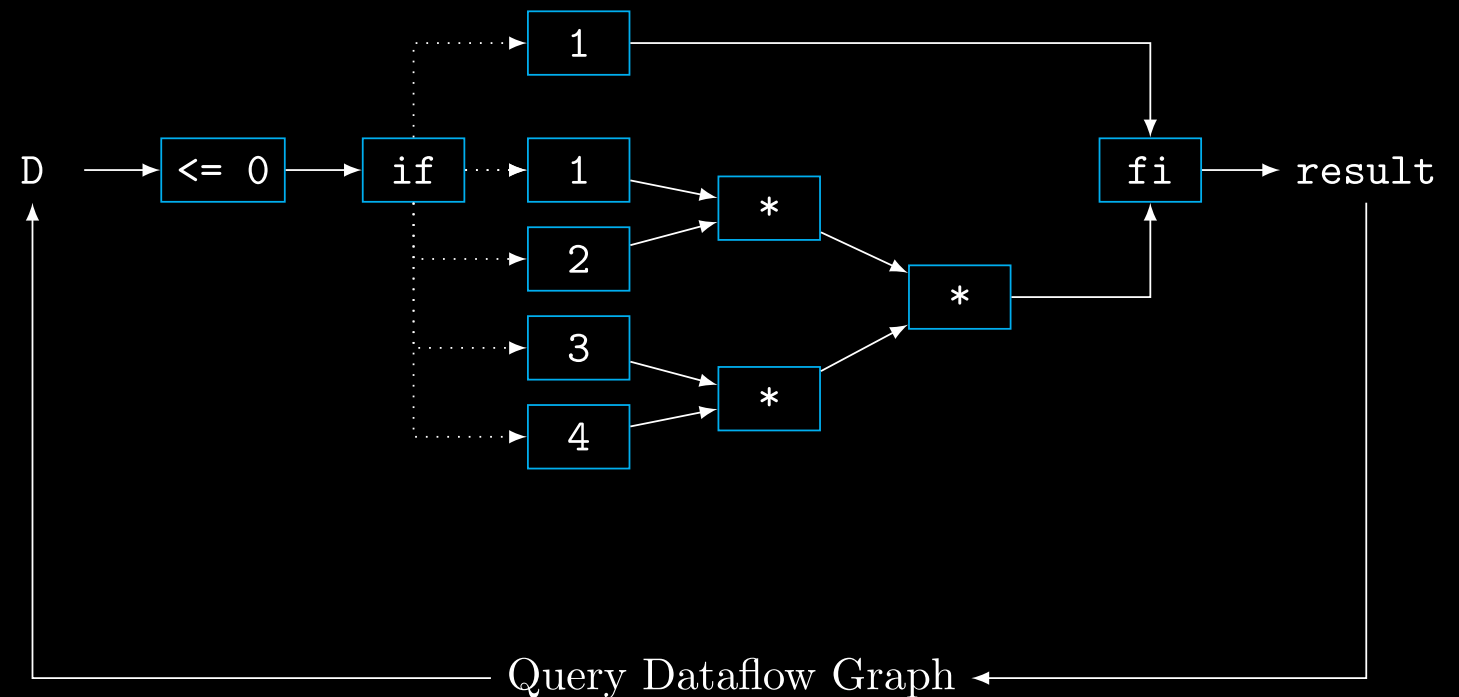
SQL / C/C++, Java, etc.



Example UDF:

```
function factorial(D) {  
  if (D <= 0) {  
    return 1;  
  } else {  
    var result = 1;  
    for (var i = 2; i <= D; i++) {  
      result = result * i;  
    }  
    return result;  
  }  
}
```

Dataflow

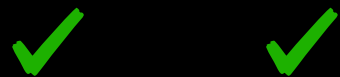


One data flow graph to rule them all!

Fully-integrated Dataflow Parallelism

Analytical Queries
(TPC-H)

SQL / JavaScript



Transactional Programs
(TPC-C)

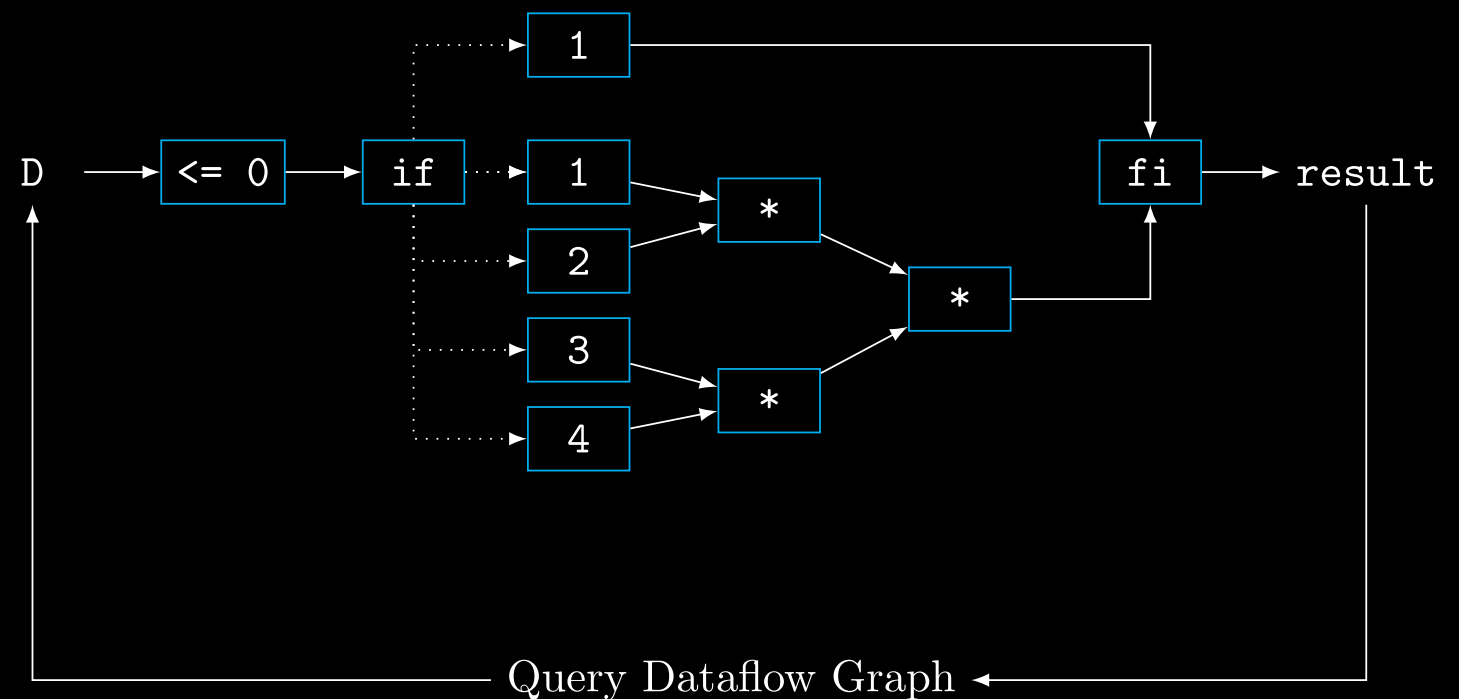
SQL / C/C++ / Java, etc.

Example UDF:

```
function factorial(D) {  
  if (D <= 0) {  
    return 1;  
  } else {  
    var result = 1;  
    for (var i = 2; i <= D; i++) {  
      result = result * i;  
    }  
    return result;  
  }  
}
```

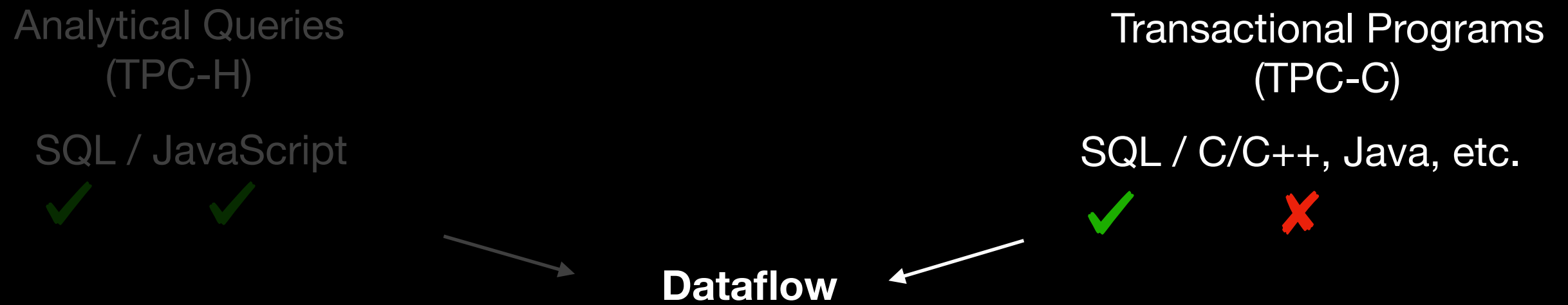
Dataflow

More optimizations apply:
pure function -> caching



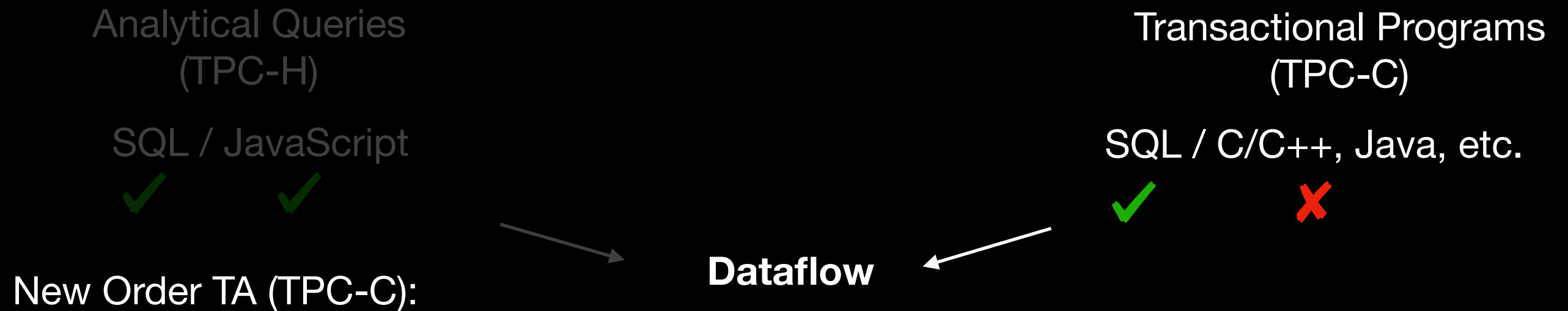
One data flow graph to rule them all!

Compiling for Efficient I/O



Challenge: Shipping computation vs. shipping data!

Compiling for Efficient I/O



Compiling for Efficient I/O

Analytical Queries
(TPC-H)

SQL / JavaScript



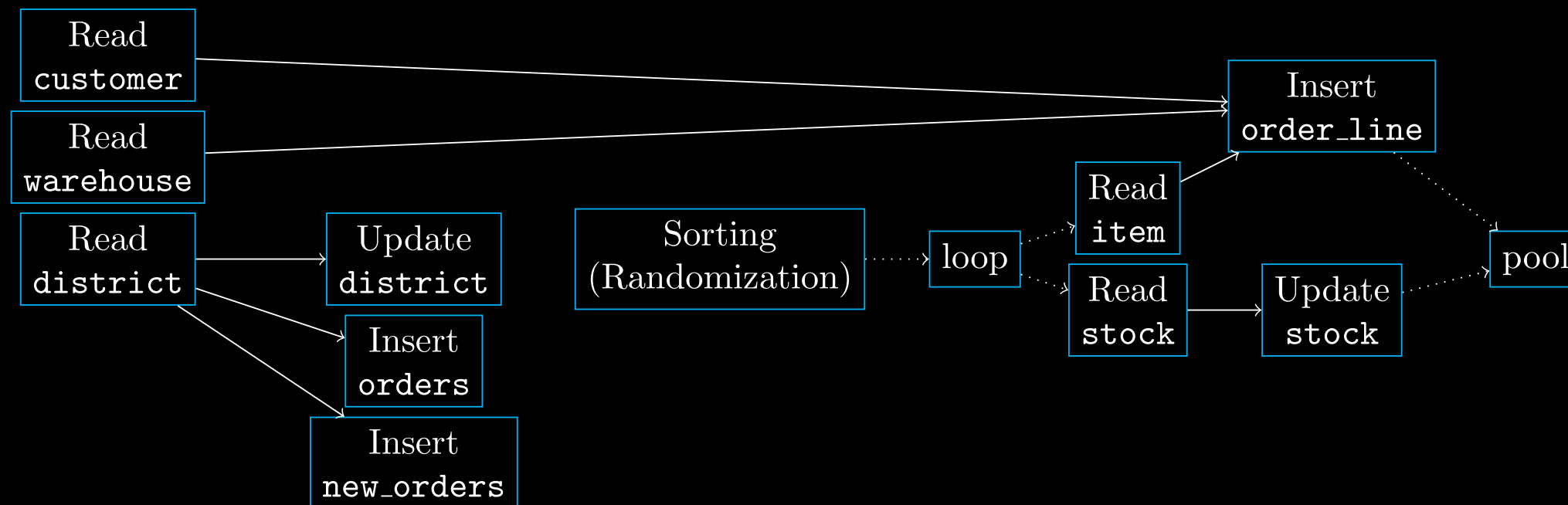
Transactional Programs
(TPC-C)

SQL / C/C++, Java, etc.



New Order TA (TPC-C):

Dataflow



Compiling for Efficient I/O

Analytical Queries
(TPC-H)

SQL / JavaScript



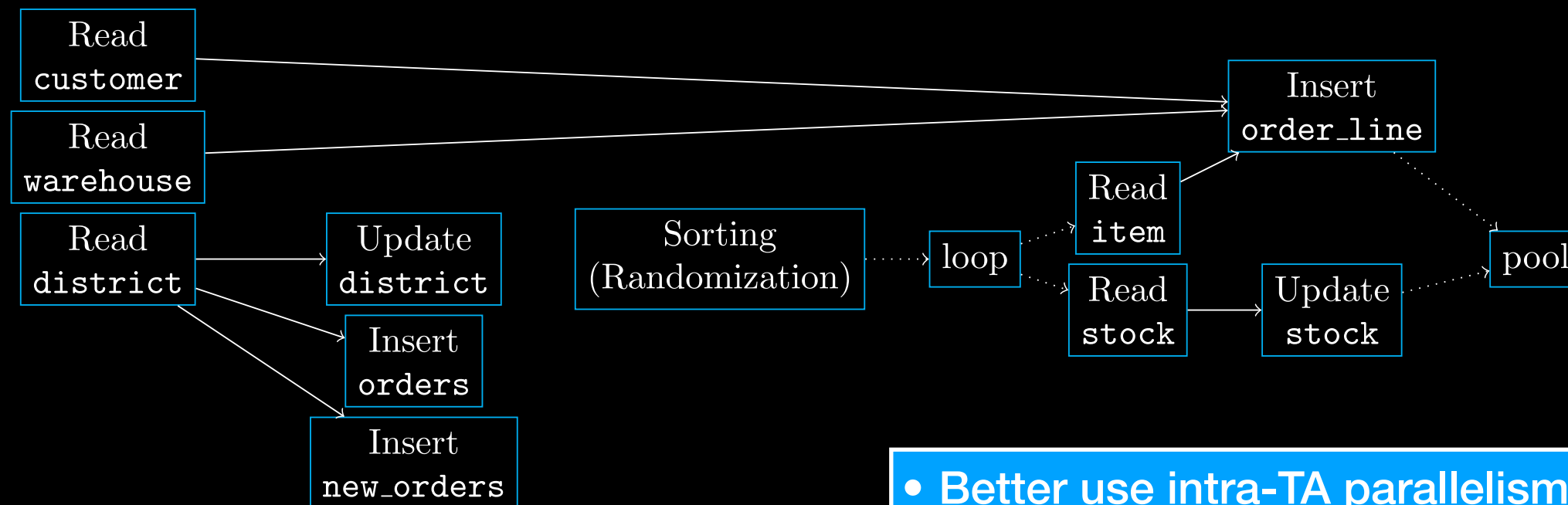
Transactional Programs
(TPC-C)

SQL / C/C++, Java, etc.



New Order TA (TPC-C):

Dataflow



- Better use intra-TA parallelism!

Compiling for Efficient I/O

Analytical Queries
(TPC-H)

SQL / JavaScript



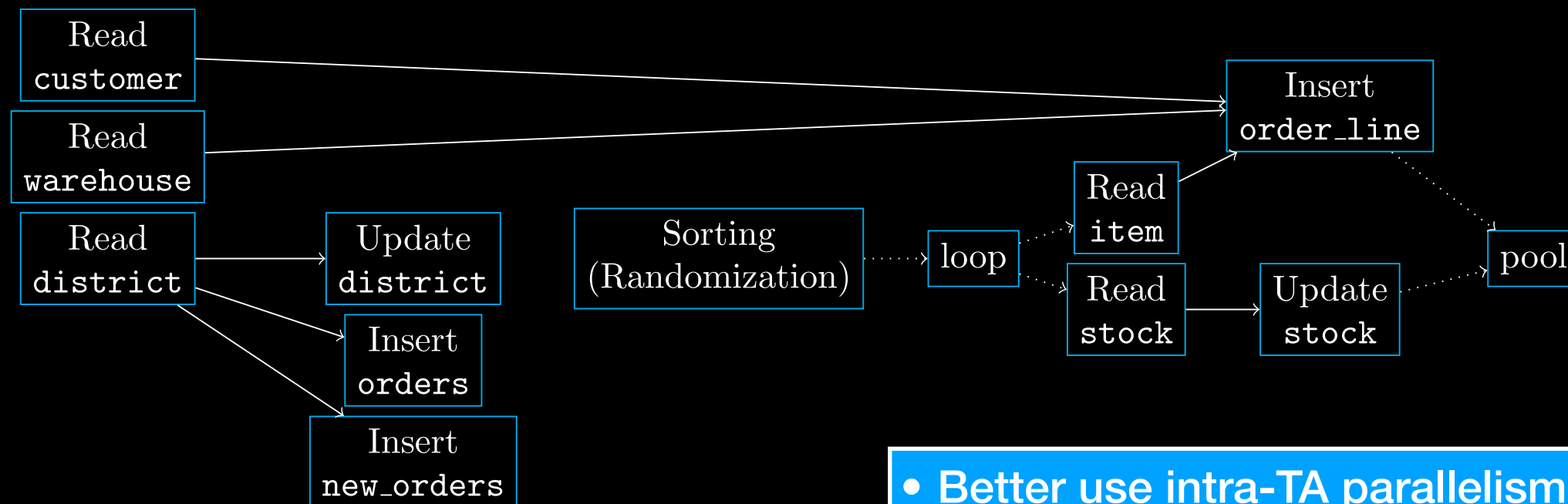
Transactional Programs
(TPC-C)

SQL / C/C++, Java, etc.



New Order TA (TPC-C):

Dataflow



- Better use intra-TA parallelism!
- Migrate not only query but computation!

Key Insight

Instead of SQL *vs.* Java, C/C++ etc. better have a
common foundation!

Thank you!