# Unified Engine for Data Processing and AI

Xiaowei Jiang
Sept, 2018

**Alibaba** Group

**Stratify Inc.**
2000 - 2002
Member of Technical Staff

**Facebook**
2010 - 2014
Software Engineer

2002

2010

2014

Now

**Microsoft**
2002 - 2010
Principal Software Engineer

**Alibaba Group**
2014 - Now
Senior Director

Alibaba Group

聚划算  淘宝网  天猫 Tmall.com  阿里云  **Alibaba** Group  支付宝 ALIPAY  CAINIAO 菜鸟  AliExpress

EB Total    PB Everyday    1T Event/Day    472M Events/sec

**472M events/s**

**Sub-second latency**

**Exactly-once**

**Highly Available**

# Flink Architecture

**Alibaba** Group

**Different API for Stream and Batch Processing**

| Table API & SQL |
| :---: |
| Relational |

| DataStream API | DataSet API |
| :---: | :---: |
| Stream Processing | Batch Processing |

| Runtime |
| :---: |
| Distributed Streaming DataFlow |

| Local | Cluster | Cloud |
| :---: | :---: | :---: |
| Single JVM | Standalone/YARN | GCE/EC2 |

## DataStream API
Stream Processing

## DataSet API
Batch Processing

Transformation

Operator Tree

StreamGraph

Batch Plan

Optimized Plan

**API**

**Very different code for Stream & Batch**

Job Graph

Stream Task & OP

Batch Task & Driver
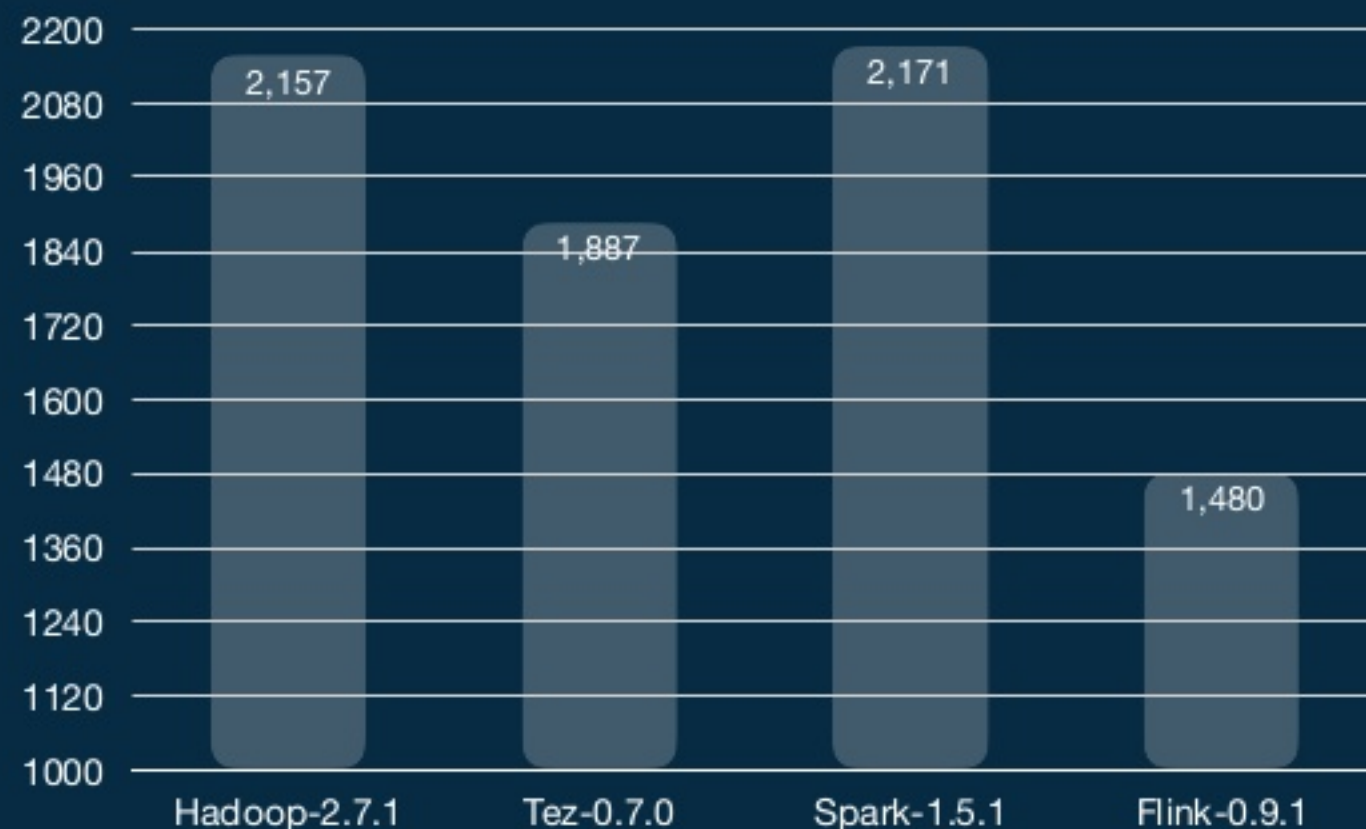
**Runtime**

Stream Processing Engines

*Apache Flink is the most sophisticated open-source Stream Processor*

Batch Processing Engines



Can Apache Flink become the
most sophisticated
open-source batch processor?

Alibaba Group

## Result of sorting 80GB/node (3.2TB)



**Flink is the fastest due to its pipelined execution**

Tez and Spark do not overlap 1st and 2nd stages
MapReduce is slow despite overlapping stages

*A Comparative Performance Evaluation of Flink, Dongwon Kim,  POSTECH, Flink Forward 2015*

Declarative

Optimizable

Understandable

Stable

Unified

SELECT h.p
FROM company
WHERE o.p
ORDER BY 2

**Alibaba** Group

Stream Join

Agg w/ Retraction

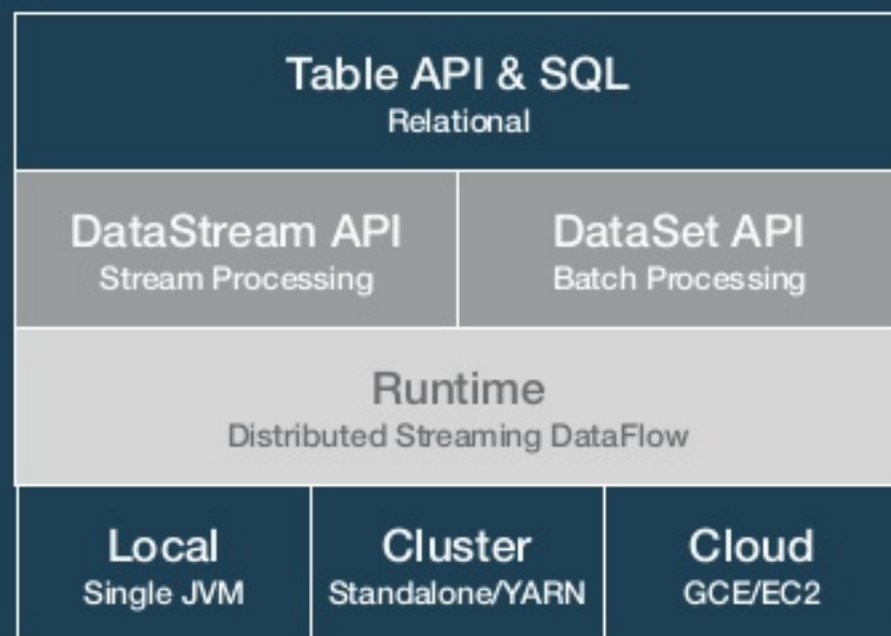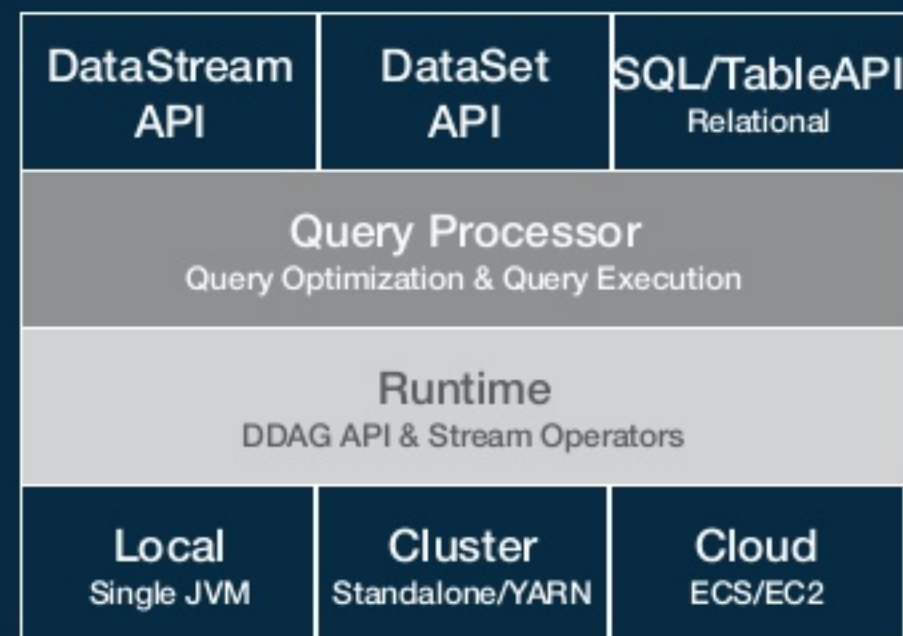Window

UDX

DDL Support

Connector

**TPCH & TPCDS**

**2K+ CASE & SQL**

## Unified Operator Abstraction

- Operators can choose inputs
- Operators can be chained easily
- Helps batch as well as streaming

Alibaba Group

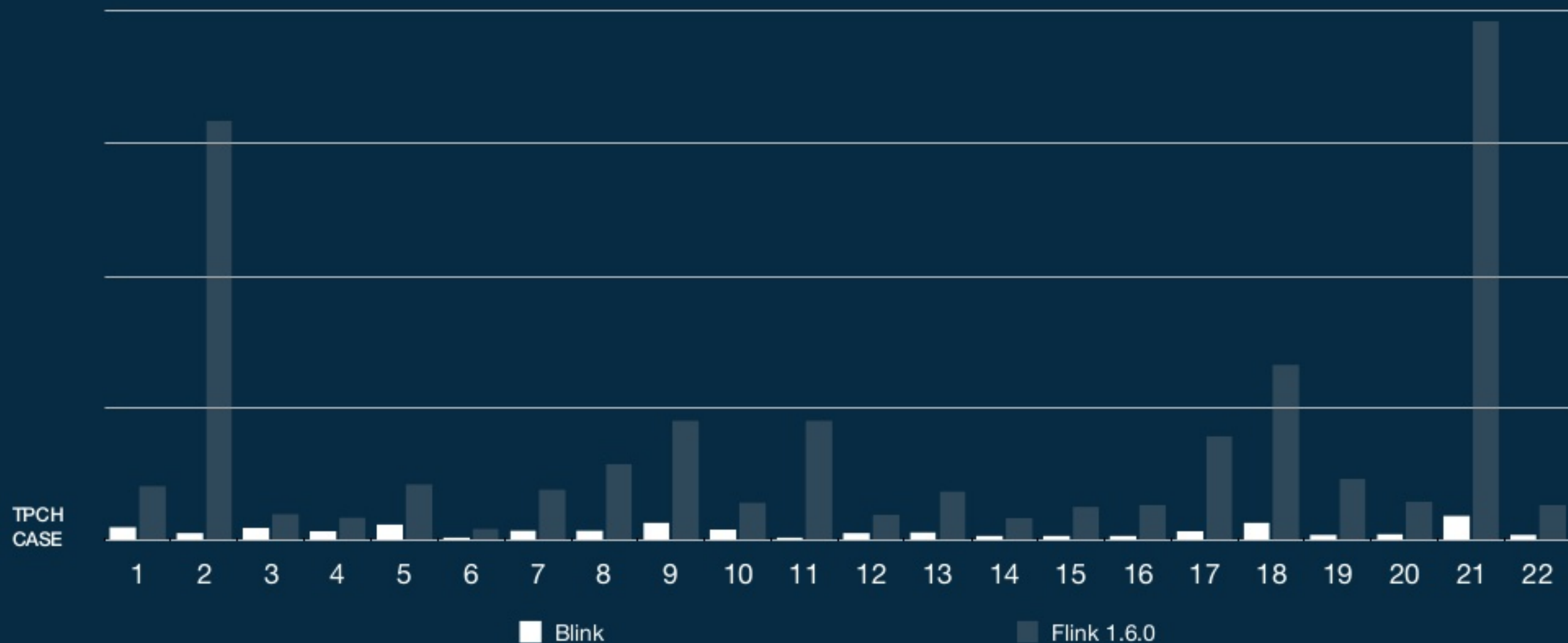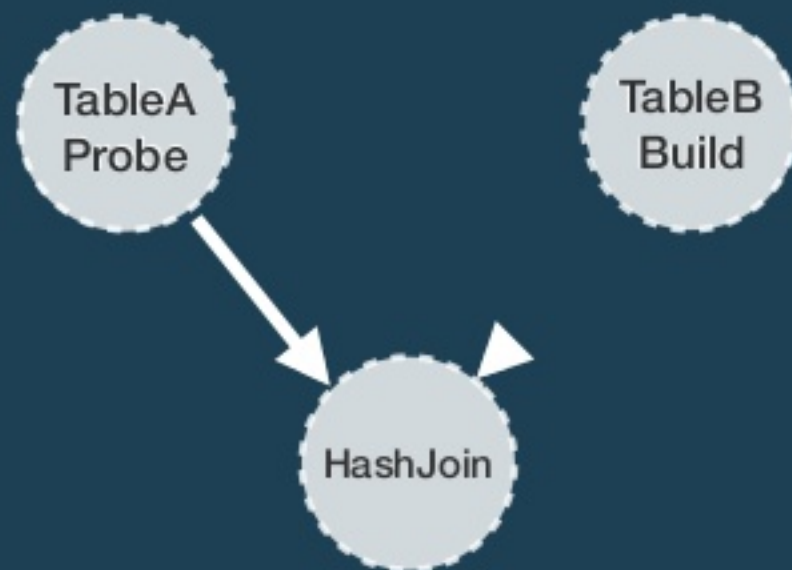| Table API & SQL | |
| --- | --- |
| Relational | |
| DataStream API | DataSet API |
| Stream Processing | Batch Processing |
| Runtime | |
| Distributed Streaming DataFlow | |
| Local | Cluster | Cloud |
| Single JVM | Standalone/YARN | GCE/EC2 |

Old Design

| DataStream API | DataSet API | SQL/TableAPI |
| --- | --- | --- |
| | | Relational |
| Query Processor | | |
| Query Optimization & Query Execution | | |
| Runtime | | |
| DDAG API & Stream Operators | | |
| Local | Cluster | Cloud |
| Single JVM | Standalone/YARN | ECS/EC2 |

New Design

# Batch Performance



**TPCH Performance** *(the Lower, the Better)*

TPCH CASE: 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22

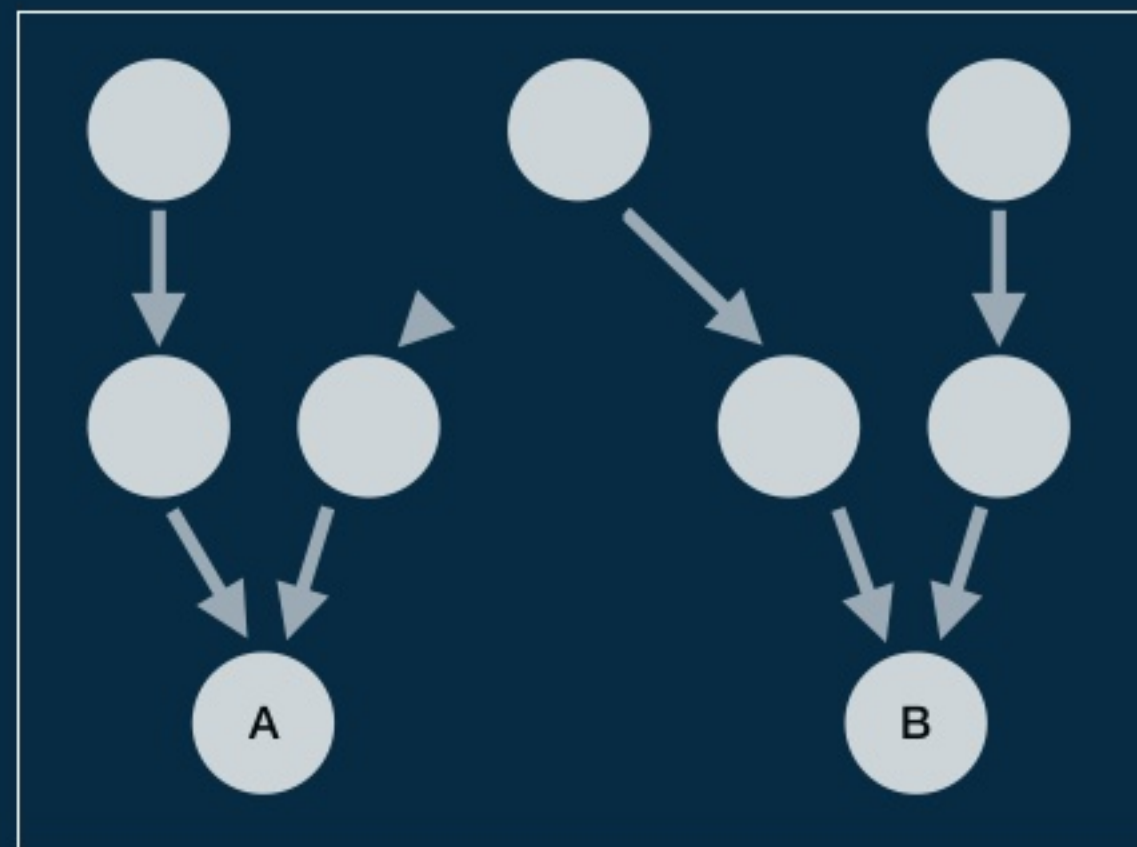■ Blink    ■ Flink 1.6.0

## Customizable Scheduling

- *Flexible control over when tasks get scheduled*
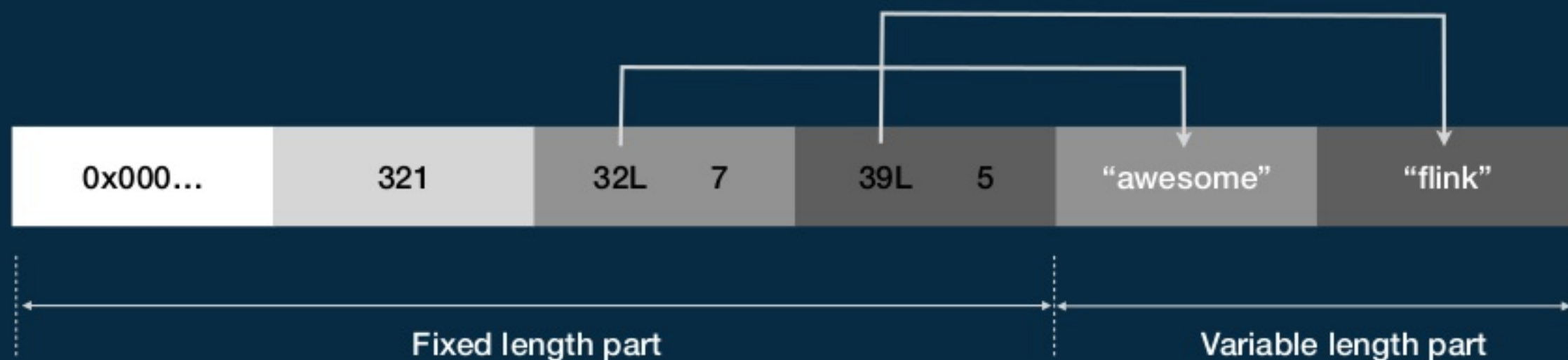- *Much better resource usage achieved*

# Flexible Chaining



Only One-Input Operators can be chained

Multi-Input Operators can also be chained

# Record Format

- Introduced new row format: BinaryRow
- Tight integration with memory management
- Avoid deserialization cost



| 0x000… | 321 | 32L 7 | 39L 5 | "awesome" | "flink" |

Fixed length part                                    Variable length part

**Alibaba** Group

**Expression Optimizations**

Operate directly on binary data
JVM intrinsics
Hot method codegen

**Performant Operators**

Operator codegen
HashAgg
Improved HashJoin
Semi/Anti join

**Resource Optimizations**

Stats based estimation
Dynamic memory allocation

**Alibaba** Group

### Cost Based

- Join order
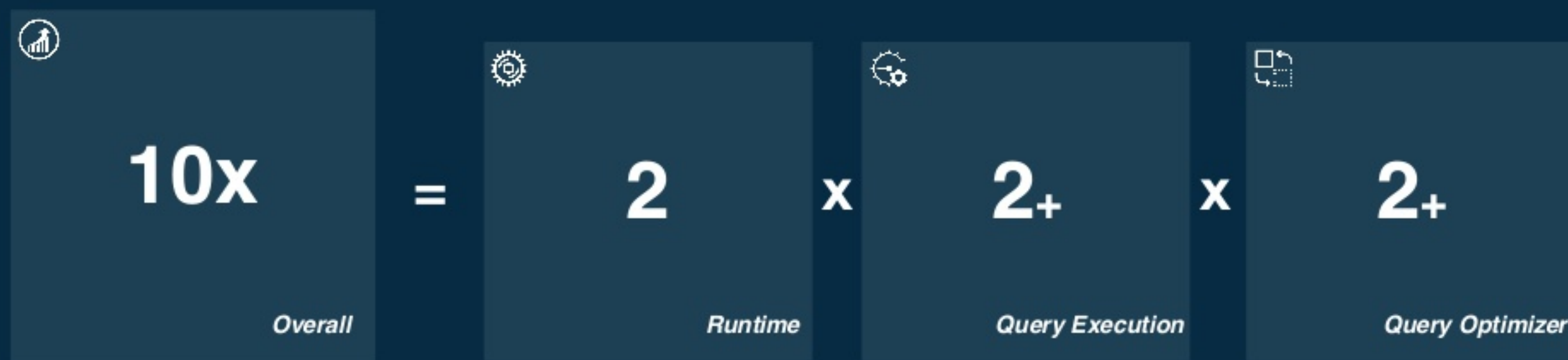- Join type
- Agg strategy
- ......

### Advanced Rules

- Subplan reuse
- Join condition expansion
- Shuffle removal
- Distinct Agg rewrite
- ....

### Rich Stats

- NDV
- NULL count
- Avg length
- Max length
- Min
- Max

Unified Engine

Functionality

Performance

**Reliability**

Alibaba Group

**Failover**

- Region Based Failover
- JM Failover
- Blacklist
- ....

**Shuffle**

- Decoupled from TM
- Yarn Shuffle Service
- Async mode

*More Details?*

Sept 4th, 2018
5:10 PM - 5:50 PM
Maschinenhaus
Feng Wang, Alibaba

**Runtime Improvements for Flink Batch Processing**

Next Steps

## Grand Unification of Data Processing

Switch between batch processing and streaming seamlessly

## Flink Machine Learning/AI

PyFlink, TableAPI, DL Integration, Flink ML Improvements

**Flink Forward China**

**Dec 20th-21st @ Beijing National Conference Center**
First Flink Forward Conference in Asia, 3000+ participants expected

**Flink Community**
Joint efforts by all major players in Flink community from China

**Call For Talks & Sponsors**
Submit your talks to flink-forward-china@list.alibaba-inc.com