# WHEN TWO CHOICES ARE NOT ENOUGH
## Balancing at Scale in Distributed Stream Processing

Muhammad Anis Uddin Nasir[1S], Gianmarco de Francisci Morales[3c], Nicolas Kourtellis[2¥], Marco Serafini[3€]

[1] KTH Royal Institute of Technology, [2] Telefonica Research Barcelona, [3] Qatar Computing Research Institute

anisu@kth.se[S], gdfm@acm.org[c], nicolas.kourtellis@telefonica.com[¥], mserafini@qf.org.qa[€]
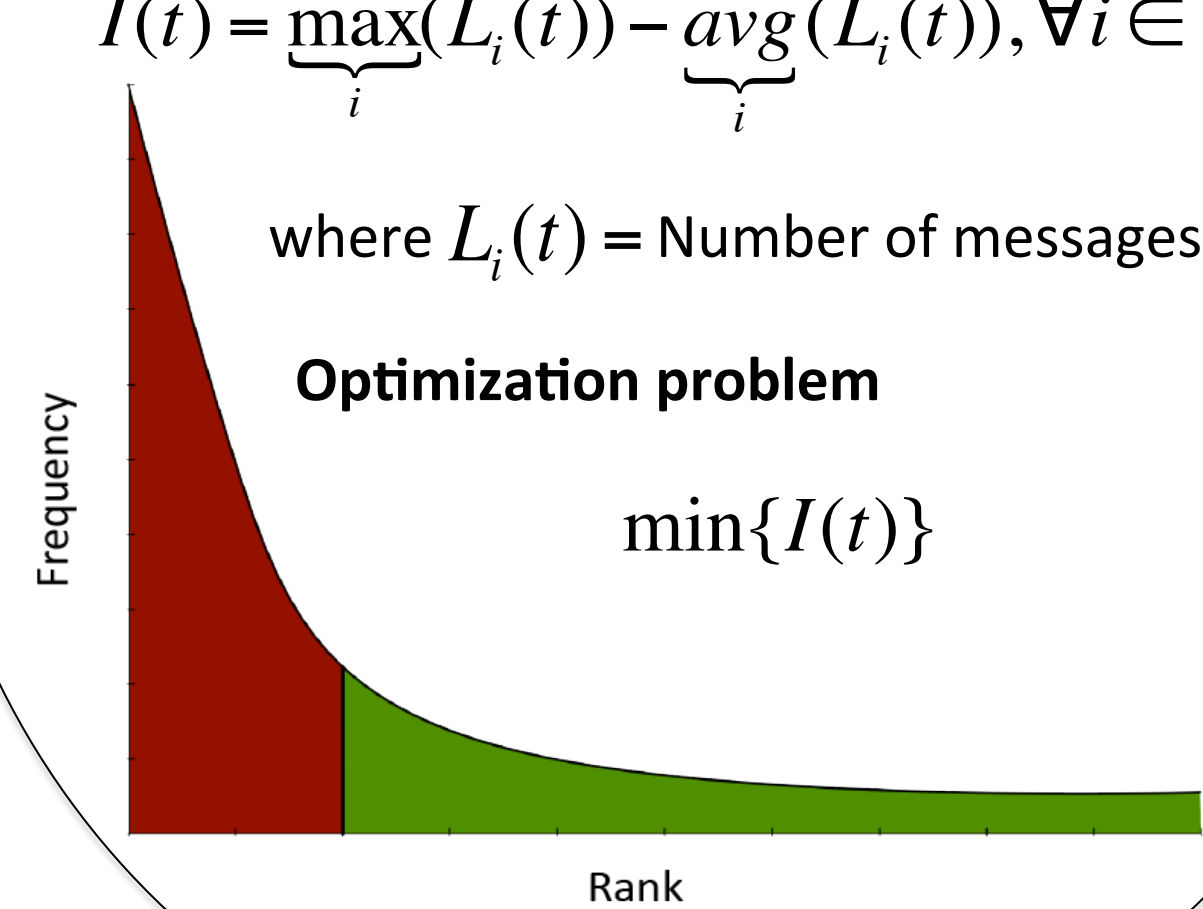
## Abstract

- Given a **highly skewed input stream**, **minimize the load imbalance** across the machines to achieve **better hardware utilization, higher throughput and lower processing latency**

- We propose two novel algorithms : a) **D-Choices** and b**) W-Choices**. Both the algorithms operate by identifying the head of the input stream and placing the head on more than two workers

## Problem Statement

$$I(t) = \underbrace{\max_i(L_i(t))}_{} - \underbrace{avg_i(L_i(t))}_{}, \forall i \in W$$

where $L_i(t)$ = Number of messages

**Optimization problem**

$$\min\{I(t)\}$$



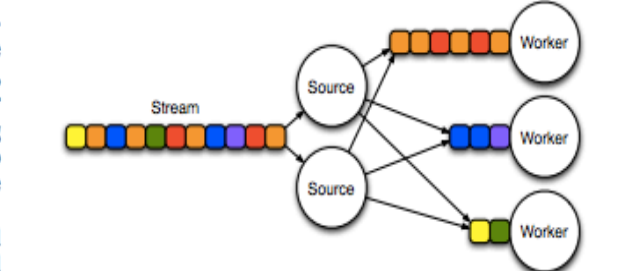Imbalance — Load — Worker (Frequency vs Rank)

## Past Work: ICDE 2015

The Power of Both Choices: Practical Load Balancing for Distributed Stream Processing Engines

Muhammad Anis Uddin Nasir[#], Gianmarco De Francisci Morales[†], David Garcia-Soriano[†]
Nicolas Kourtellis[†], Marco Serafini[§]

[#] KTH Royal Institute of Technology, Stockholm, Sweden
[†] Yahoo Labs, Barcelona, Spain
[§] Qatar Computing Research Institute, Doha, Qatar

Fig. 1: Load imbalance generated by skew in the key distribution when using key grouping. The color of each message represents its key.
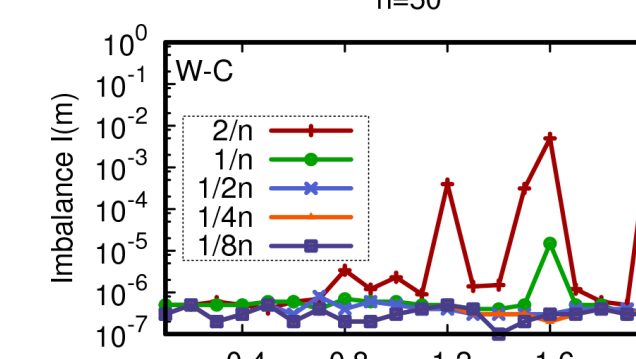
## Solution Overview



splits head from the tail — Head / Tail — k1 (p1) / Rank / kn (pn)

## Solution Overview

**Observation**

PKG guarantees nearly perfect load balance for **p1 ≤ 1/(5n)**

**Solution**

Stream Summary to handle the Head



**D-Choices:**
- adapts to the frequencies
- allows subset of workers

**W-Choices:**
- Independent of frequencies
- allows all the workers

**Algorithm:**
- Assign head to the set of d/w workers
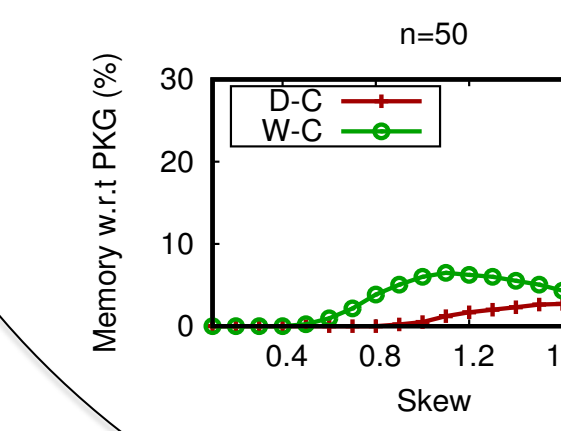- Handle tail using power of both choices

## Solution Overview

**Pseudocode**

```
Algorithm 1: Stream partitioning algorithm.
upon message m = ⟨k, v⟩
    H ← UpdateSpaceSaving(k)
    d ← 2 // Default as in PKG
    if k ∈ H then
        if D-CHOICES then
            d ← FindOptimalChoices()
        else if W-CHOICES then
            d ← n
    w ← MinLoad(F_1(k), ..., F_d(k))
    send(w, m)
```

**Memory Overhead**



## Experimental Setup

**Dataset**

| Dataset | Symbol | Messages | Keys | $p_1(\%)$ |
|---------|--------|----------|------|-----------|
| Wikipedia | WP | 22M | 2.9M | 9.32 |
| Twitter | TW | 1.2G | 31M | 2.67 |
| Cashtags | CT | 690k | 2.9k | 3.29 |
| Zipf | ZF | $10^7$ | $10^4, 10^5, 10^6$ | $\propto \frac{1}{\sum \frac{1}{x^{-z}}}$ |

**Algorithms**

| Symbol | Algorithm | Head vs. Tail |
|--------|-----------|---------------|
| D-C | D-Choices | |
| W-C | W-Choices | Specialized on head |
| RR | Round-Robin | |
| PKG | Partial Key Grouping | Treats all keys equally |
| SG | Shuffle Grouping | |

## Experimental Results

**Estimation**



**Load Imbalance**



## Experimental Results

**Latency**



**Throughput**



**Imbalance over time**



## Analysis

$$\begin{aligned} \text{minimize}_d \quad & f(d; \mathcal{D}, \theta) = d \times |\mathcal{H}_{\mathcal{D}, \theta}| \\ \text{subject to} \quad & \mathbb{E}_d[I(m)] \leq \epsilon. \end{aligned}$$

Expanding the constraint we get:

$$\sum_{i \leq h} p_i + \left(\frac{b}{n}\right)^d \sum_{h < i \leq |H|} p_i + \left(\frac{b}{n}\right)^2 \sum_{i > |H|} p_i \leq \left(\frac{b}{n}\right) + \epsilon$$

where

$$b = n - n\left(\frac{n-1}{n}\right)^{h \times d}$$

## Conclusion

- We propose two algorithms to achieve load balance at scale for DSPEs

- Use heavy hitters to separate the head of the distribution and process on larger set of workers

- Improvement translates into 150% gain in throughput and 60% gain in latency over PKG

## Future Work

- Dynamic Load Balancing for stateful operators
  - Key Migration
  - Partition Migration
  - Queuing Theory

- Load Balancing in Heterogeneous Cluster
  - Load Prediction
  - Worker Communication