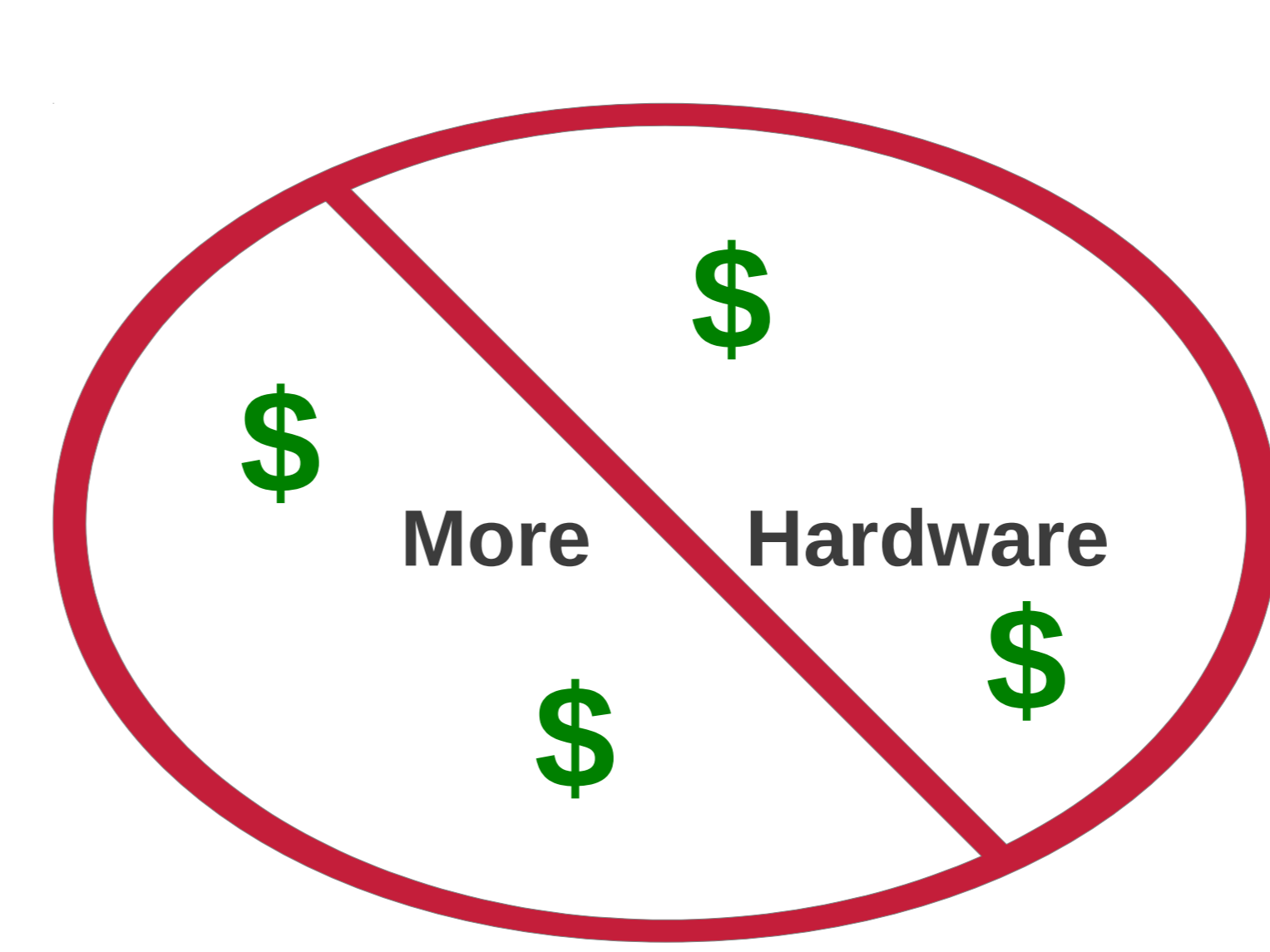
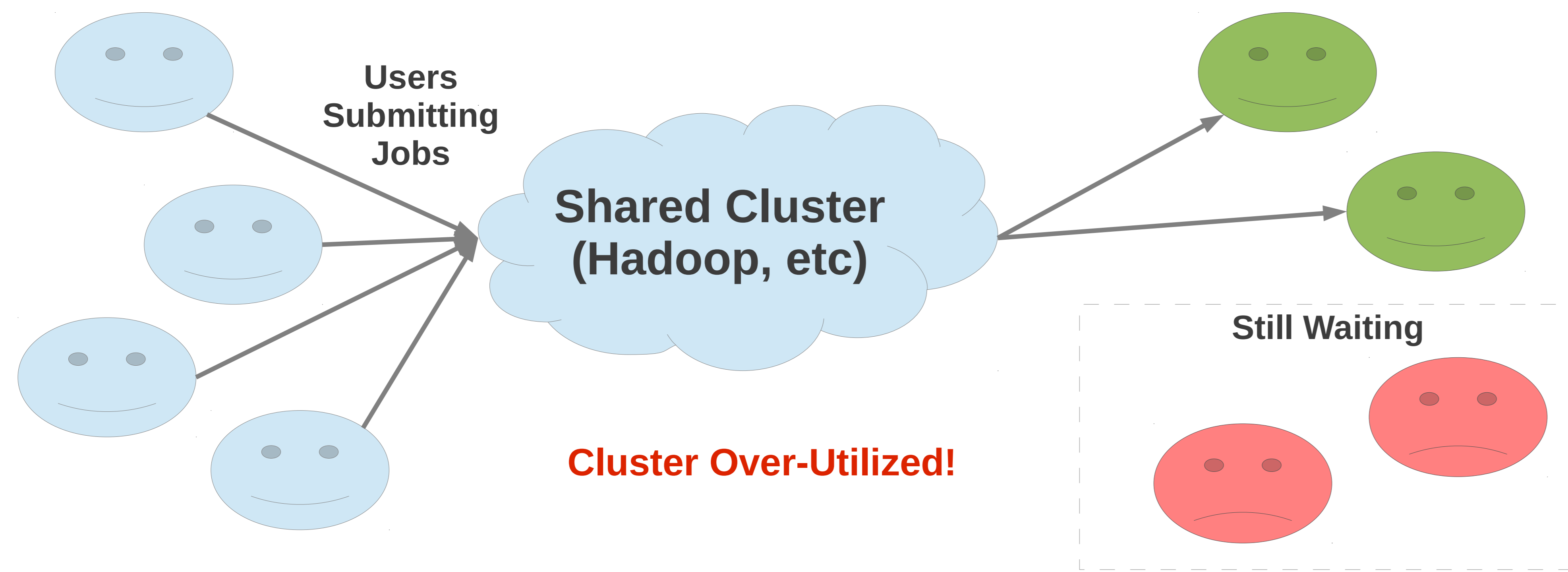
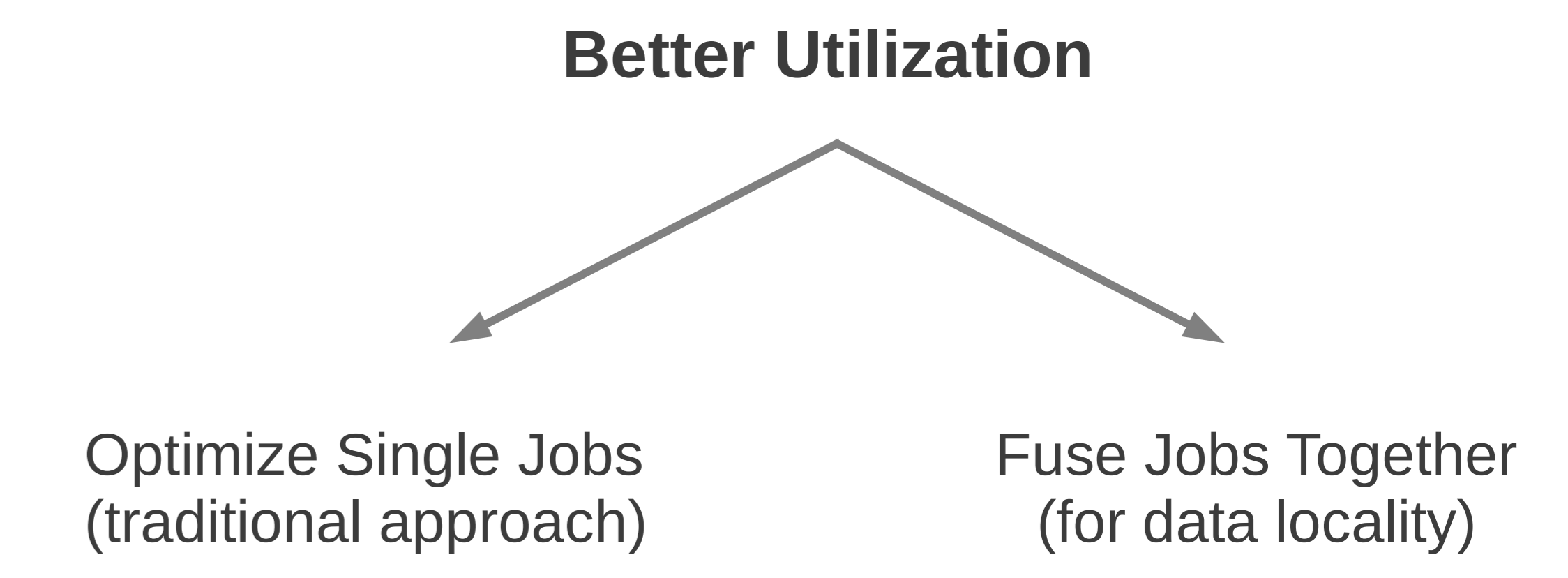


# Task Fusion: Improving Utilization of Multi-user Clusters

Motivation



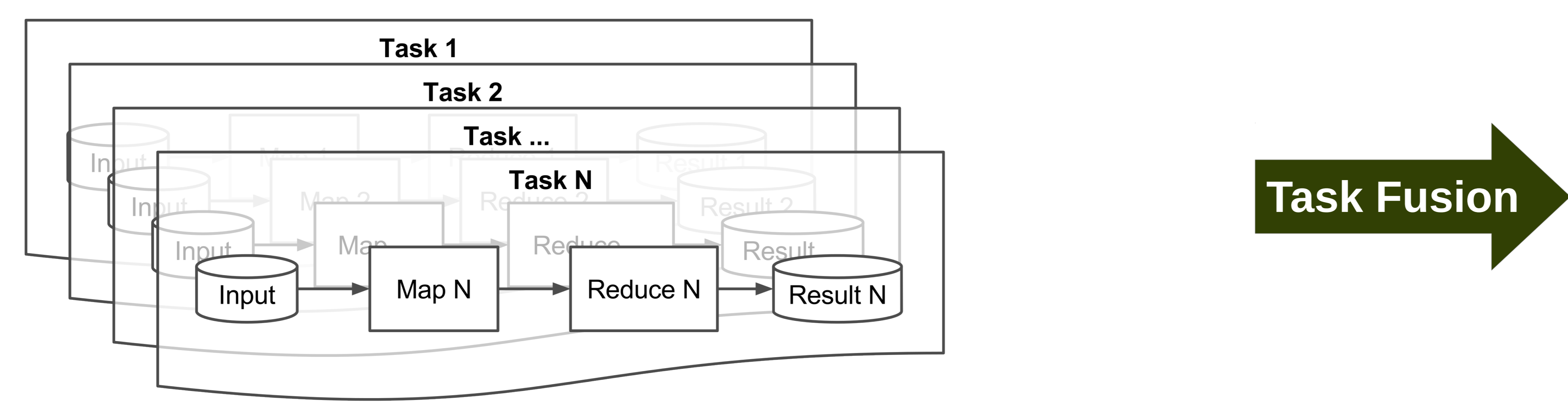
Solutions?



Insight: Load data once, run multiple analyses

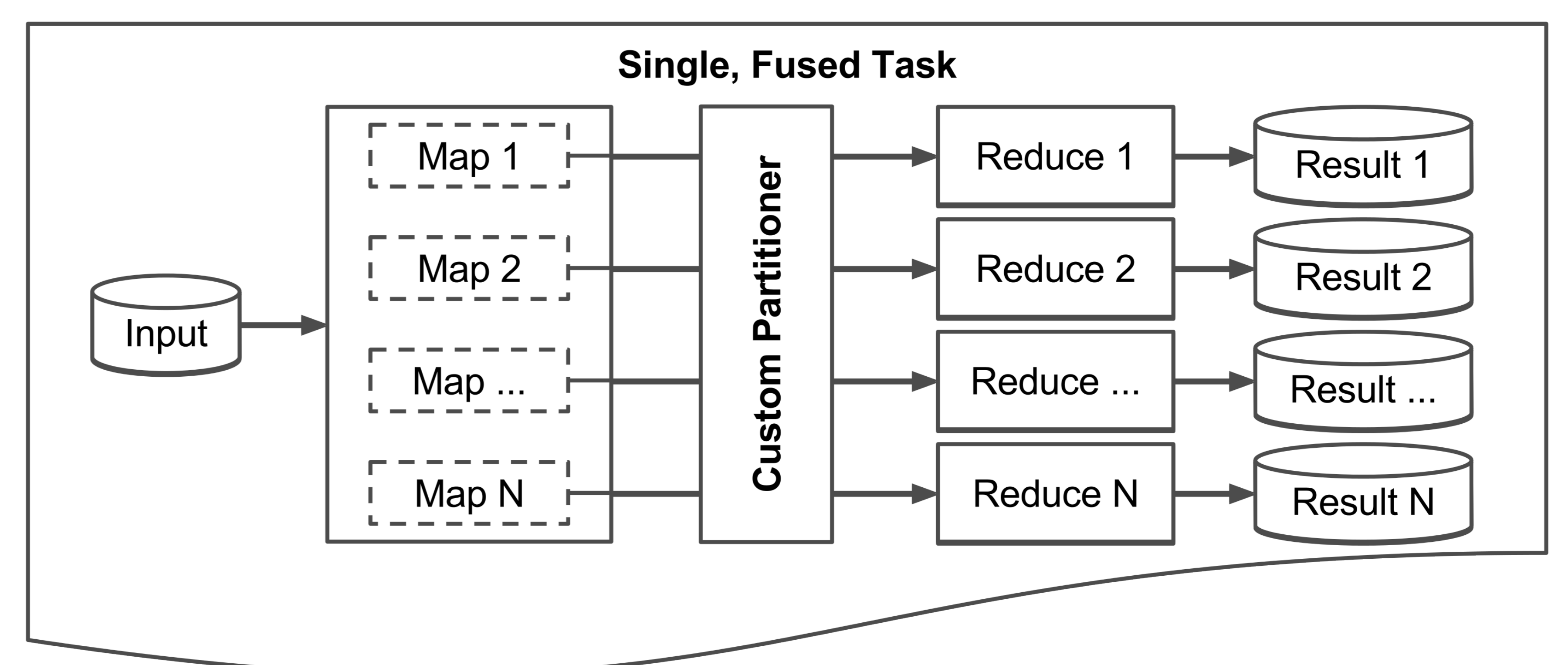
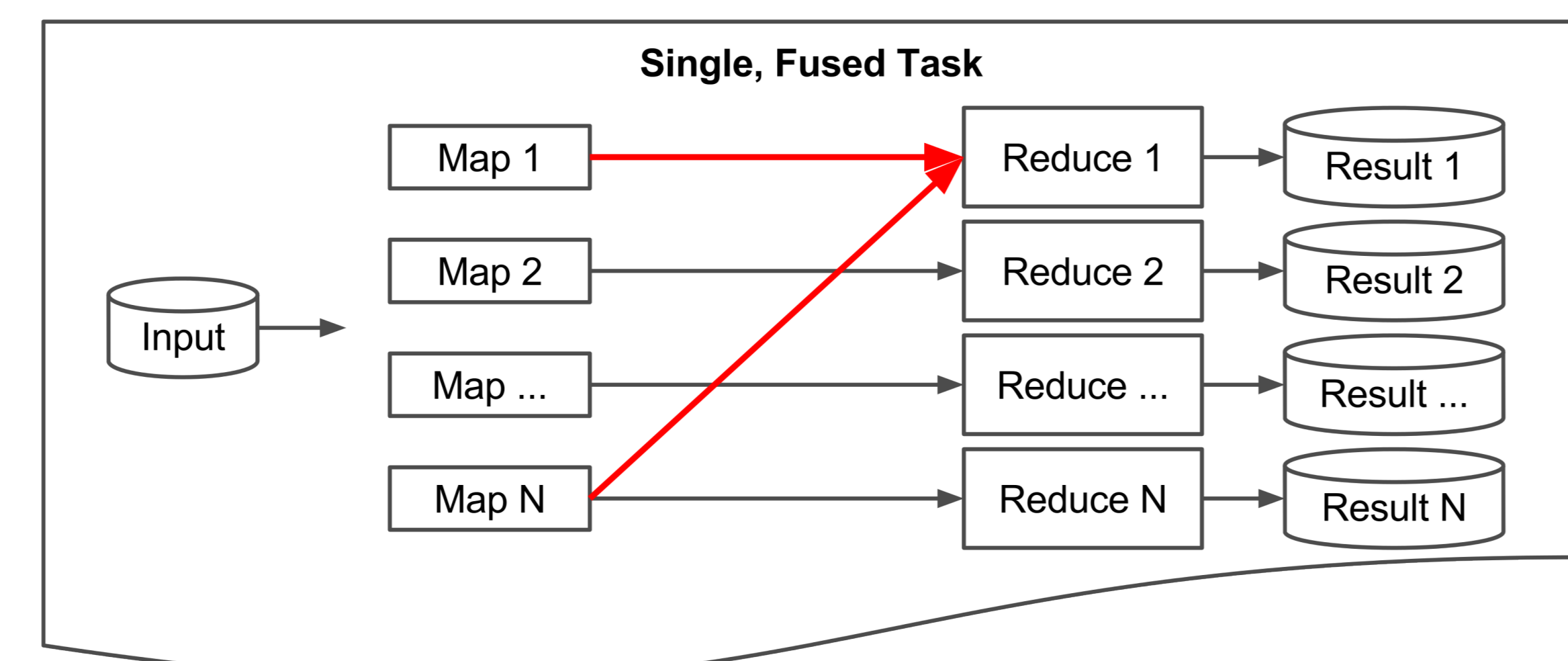
Can we automatically merge related tasks from different users?

Proposed Solution



Naive Approach

Maps from different tasks might output same keys, sending to wrong reducer



Maps modified to output tuples of (mapID, key) as keys  
Custom partitioner ensures map outputs go to correct reducer

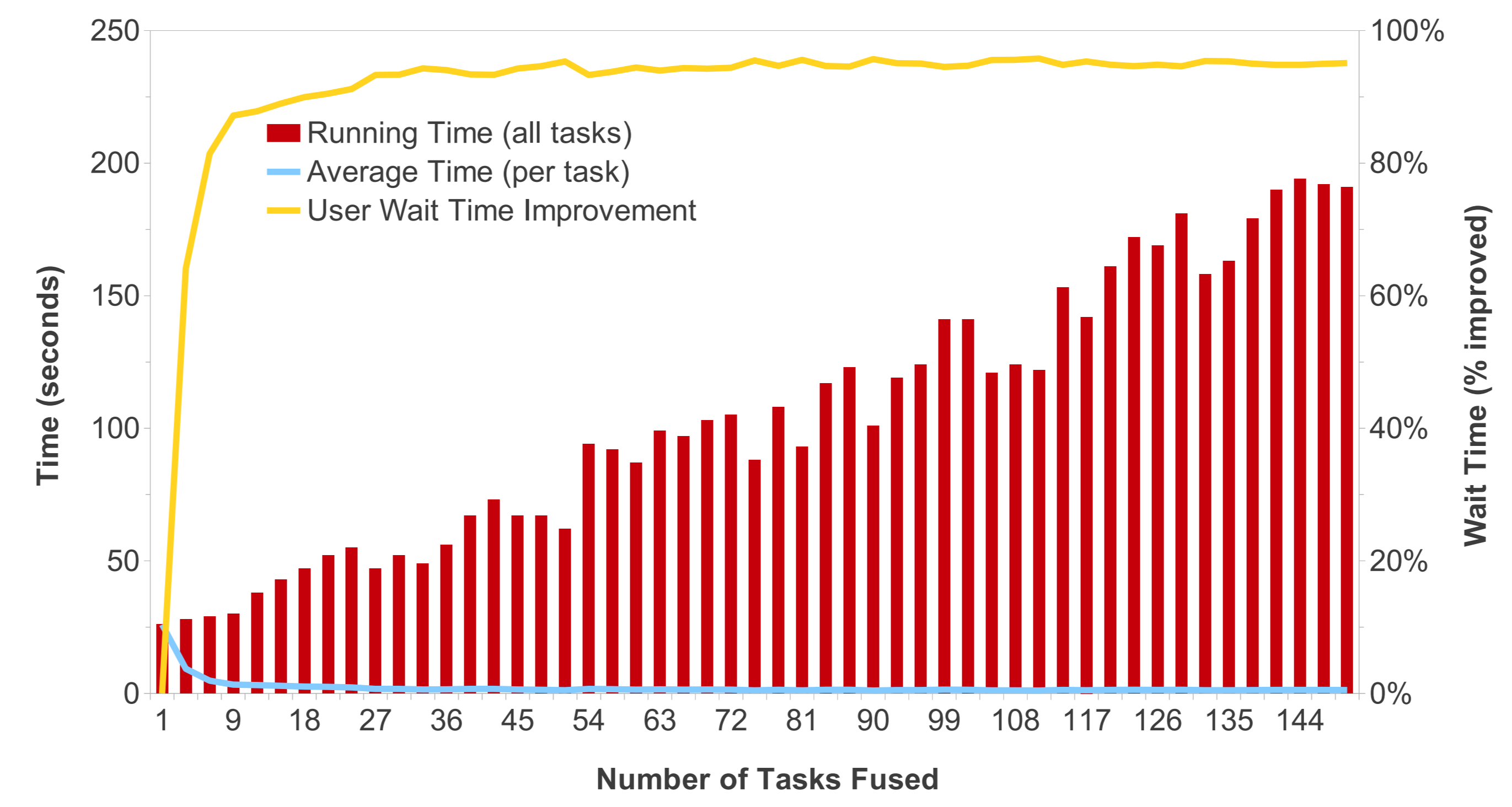
Early Results

Proof of concept implementation in Boa - <http://boa.cs.iastate.edu/>

Does task fusion increase throughput?

Task Size	# of Tasks	Times		Speedup
		No Task Fusion	Task Fusion	
Small	21	8.1m	0.8m	10.8X
Medium	22	2.3h	1.8h	1.3X
Large	18	4.6h	3.9h	1.2X
Mixed	9	1.3h	0.9h	1.4X

Does task fusion decrease user wait times?



When can I use task fusion?

Task fusion currently has the following assumptions:

1. No side-effects.
2. No shared state.
3. No dependency conflicts.

In the future...

Relax these assumptions:

- Automated program transformations
- Separate class spaces (a la OSGi)

Future Work