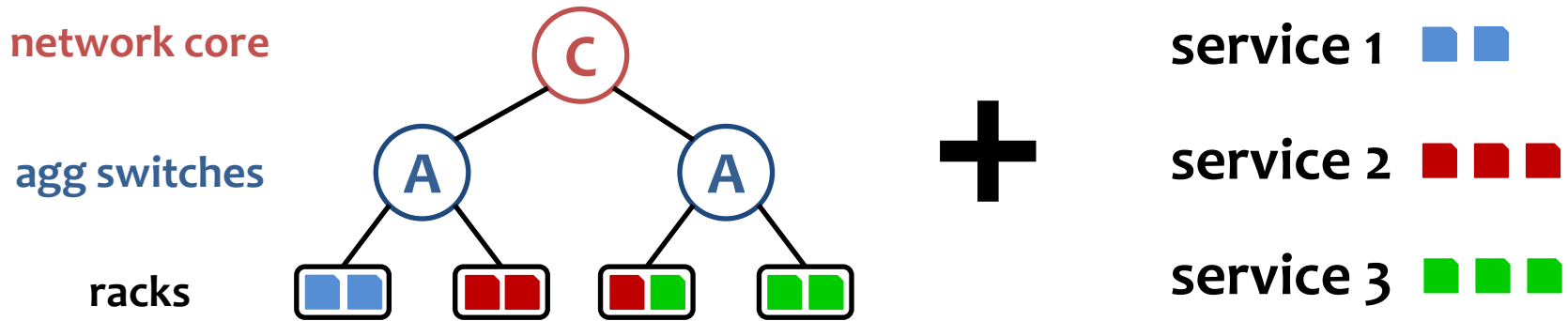


# Surviving Failures in Bandwidth-Constrained Datacenters

Peter Bodík<sup>2</sup>, Ishai Menache<sup>2</sup>, Mosharaf Chowdhury<sup>3</sup>,  
Pradeepkumar Mani<sup>1</sup>, Dave Maltz<sup>1</sup>, Ion Stoica<sup>3</sup>

Microsoft<sup>1</sup> Research<sup>2</sup>, UC Berkeley<sup>3</sup>

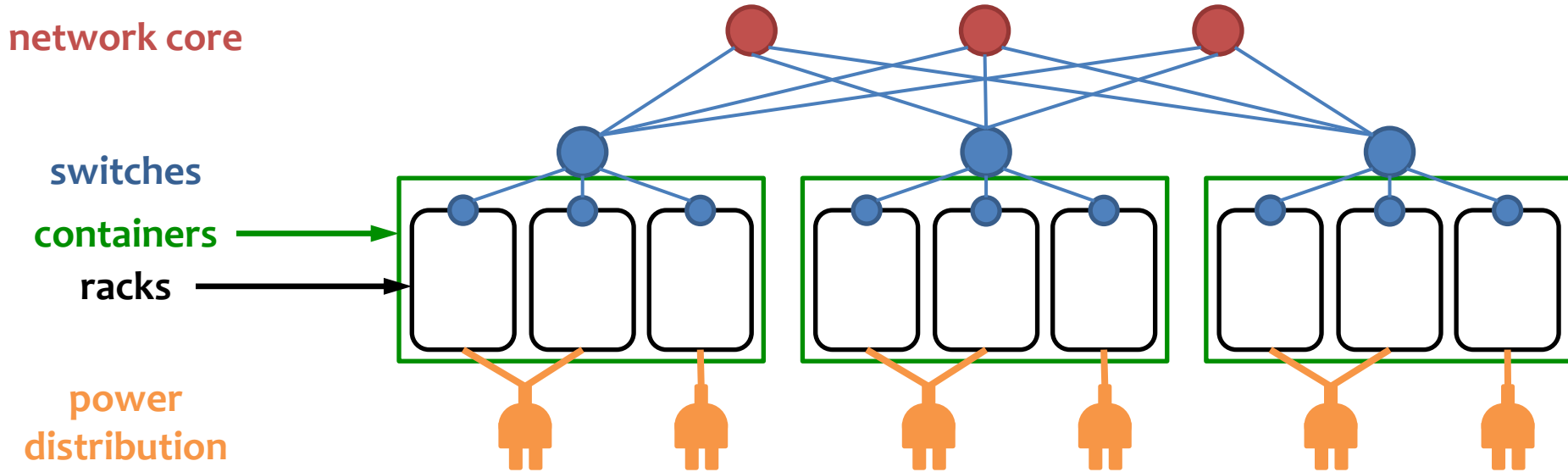
# How to allocate services to physical machines?



## Three important metrics considered together

- FT: service fault tolerance
- BW: bandwidth usage
- #M: # machine moves to reach target allocation

# FT: Improving fault tolerance of software services

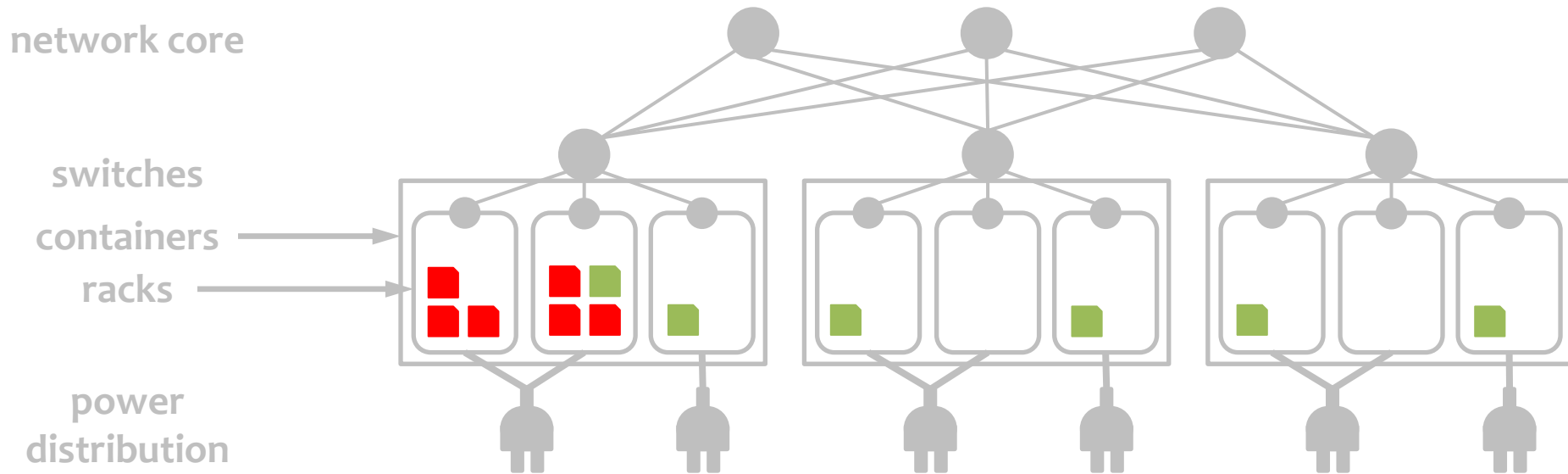


Complex fault domains: networking, power, cooling

*Worst-case survival* = fraction of service available during single worst-case failure

- corresponds to service throughput during failure

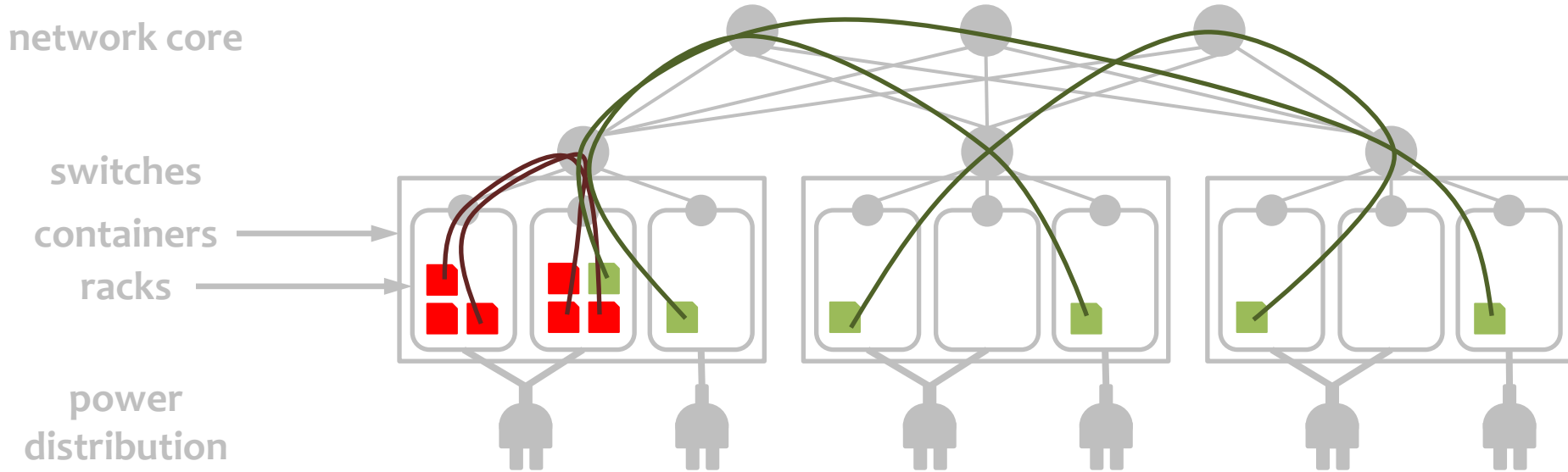
# FT: Service allocation impacts worst-case survival



## Worst-case survival:

- red service: 0% -- same container, power
- green service: 67% -- different containers, power

# BW: Reduce bandwidth usage on constrained links

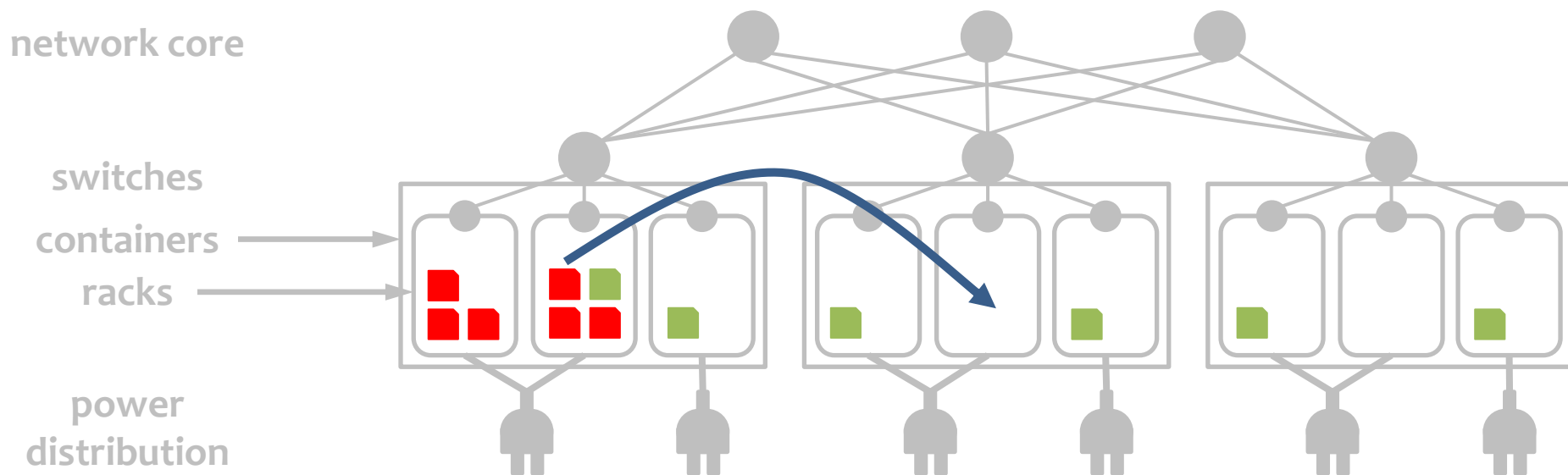


BW = bandwidth usage in the core

## Goal

- reduce cost of infrastructure
- consider other service location constraints

# #M: Need incremental allocation algorithms



## High cost of machine move

- need to deploy potentially TB of data
- warm up caches
- could take tens of minutes, impact network

# Outline

Why is it difficult?

Traffic analysis

Optimization framework

- FT + #M

- FT + BW + #M

Evaluation

# Trade-off between bandwidth usage and fault-tolerance

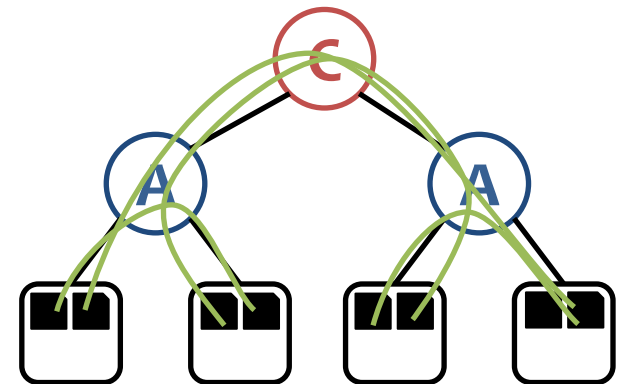
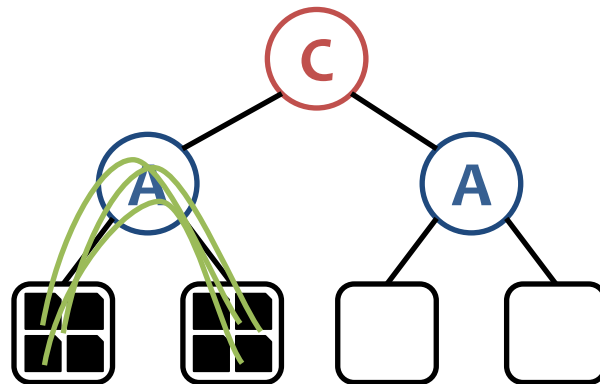
optimize for bandwidth

optimize for fault tolerance

network core

agg switches

racks



BW: utilization in core

LOW ✓

HIGH ✗

FT: fault tolerance (for agg switches)

LOW ✗

HIGH ✓

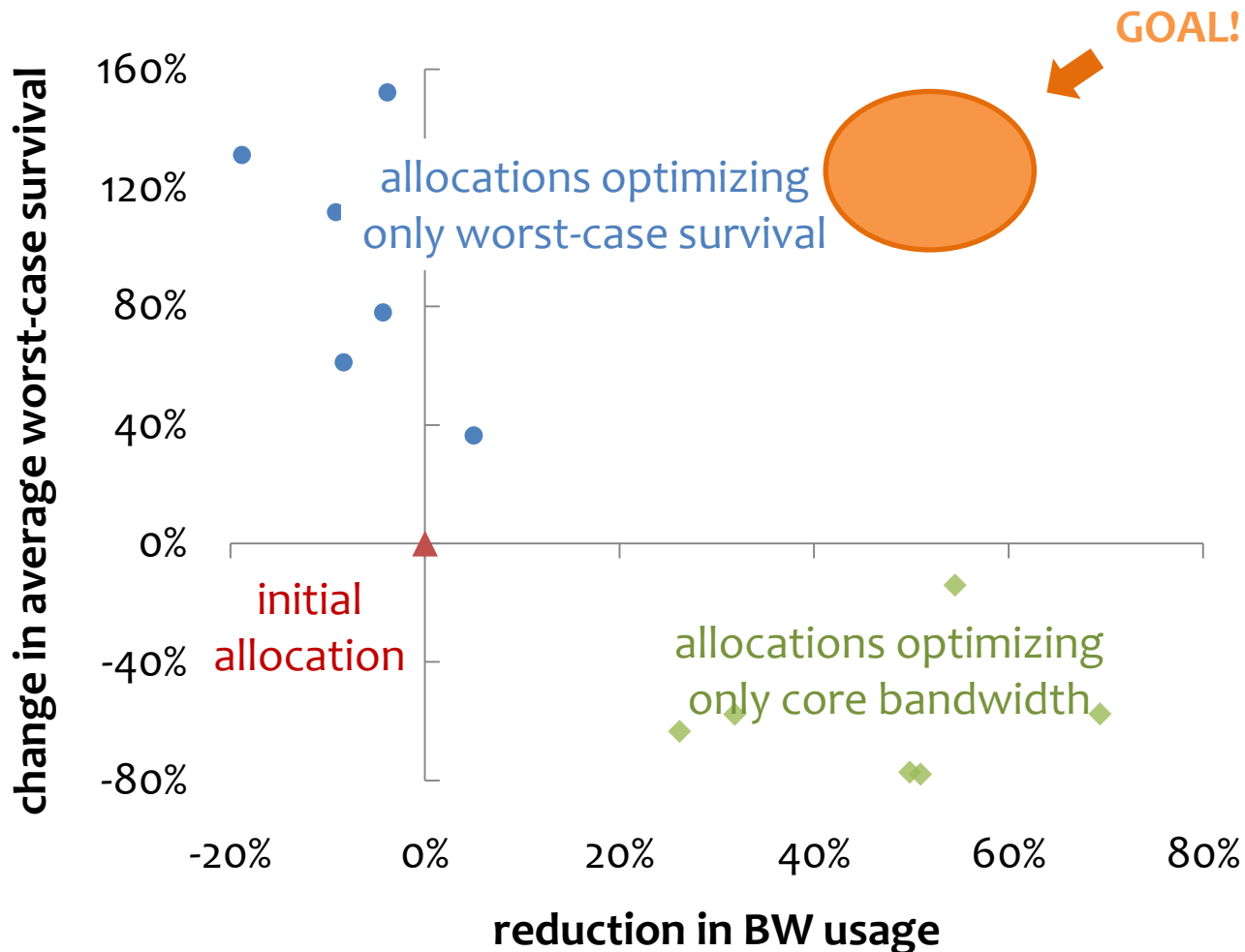
worst-case survival

0 ✗

0.5 ✓



# Optimizing for one metric degrades the other



Results from 6 Microsoft datacenters

# FT-only and BW-only are both NP-hard, hard to approximate

FT reduces to max independent set

BW reduces to min-cut in a graph

– considered previously in [Meng et al., INFOCOM'10]

Most algorithms not incremental, ignore #M

# Key insights

Improve FT using convex optimization

- local optimization leads to good solutions

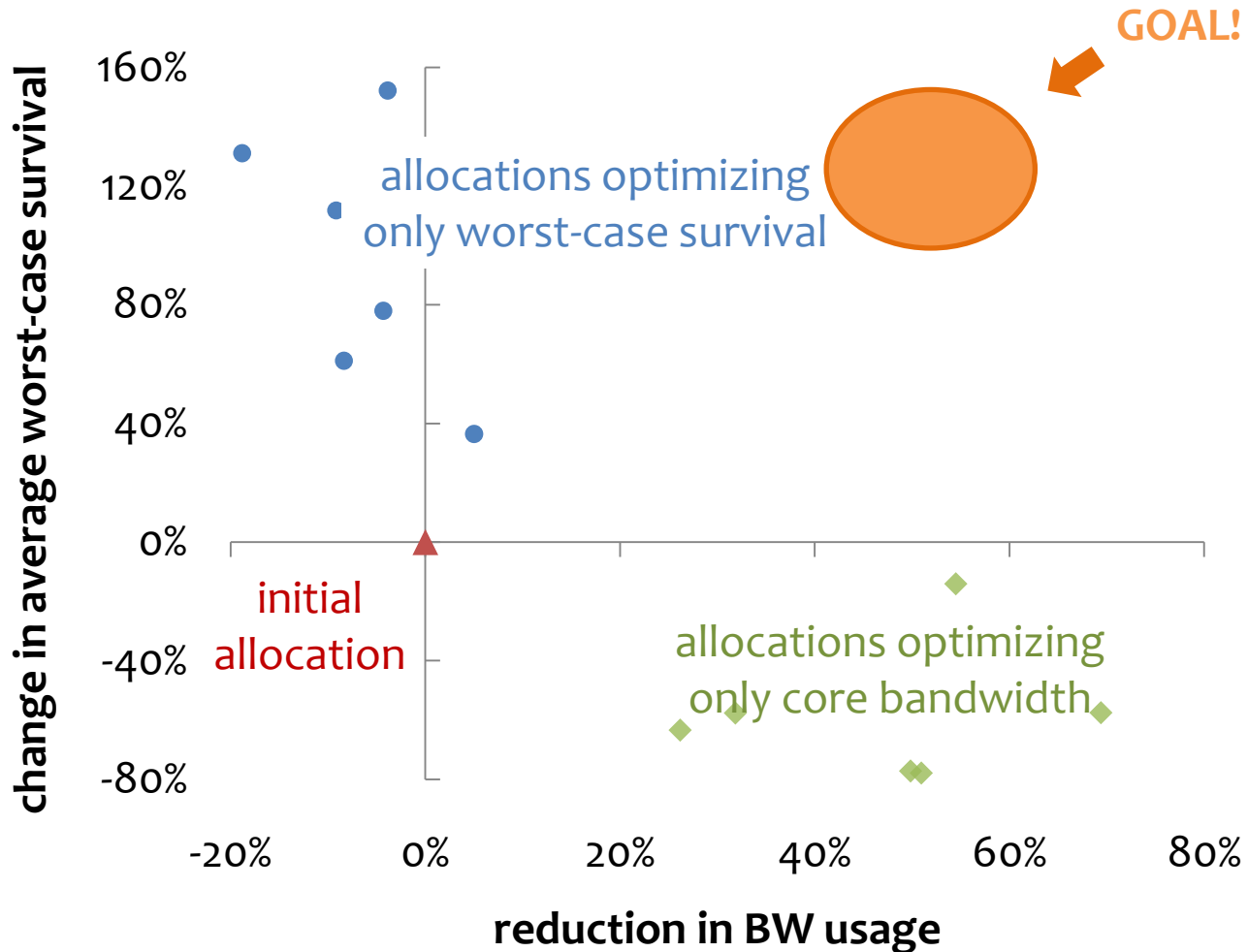
Symmetry in the optimization space

- machines, racks, containers are interchangeable

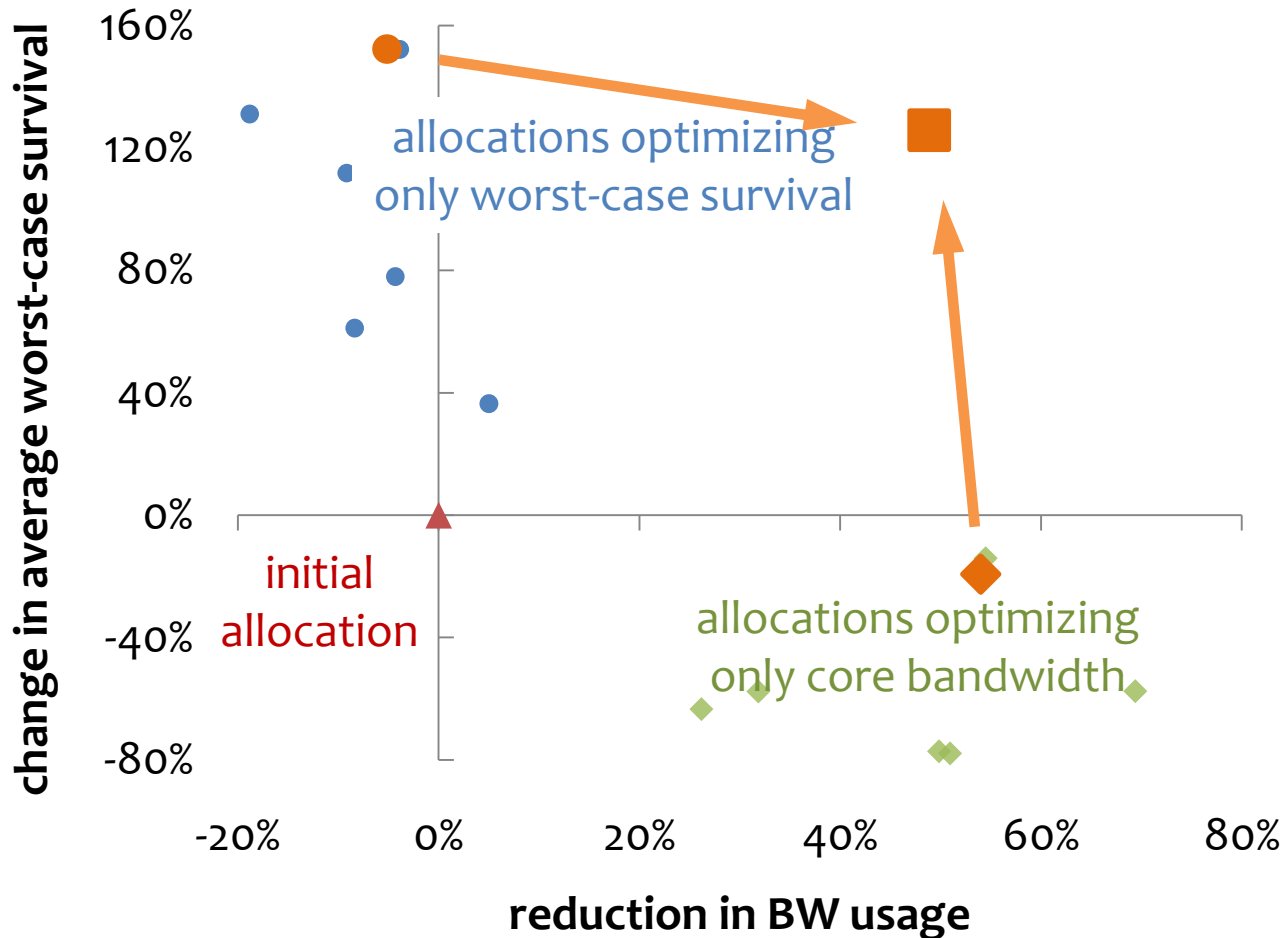
Communication pattern is very skewed

- can spread low-talkers without affecting BW

# Results preview



# Results preview



# Outline

Why is it difficult?

Traffic analysis

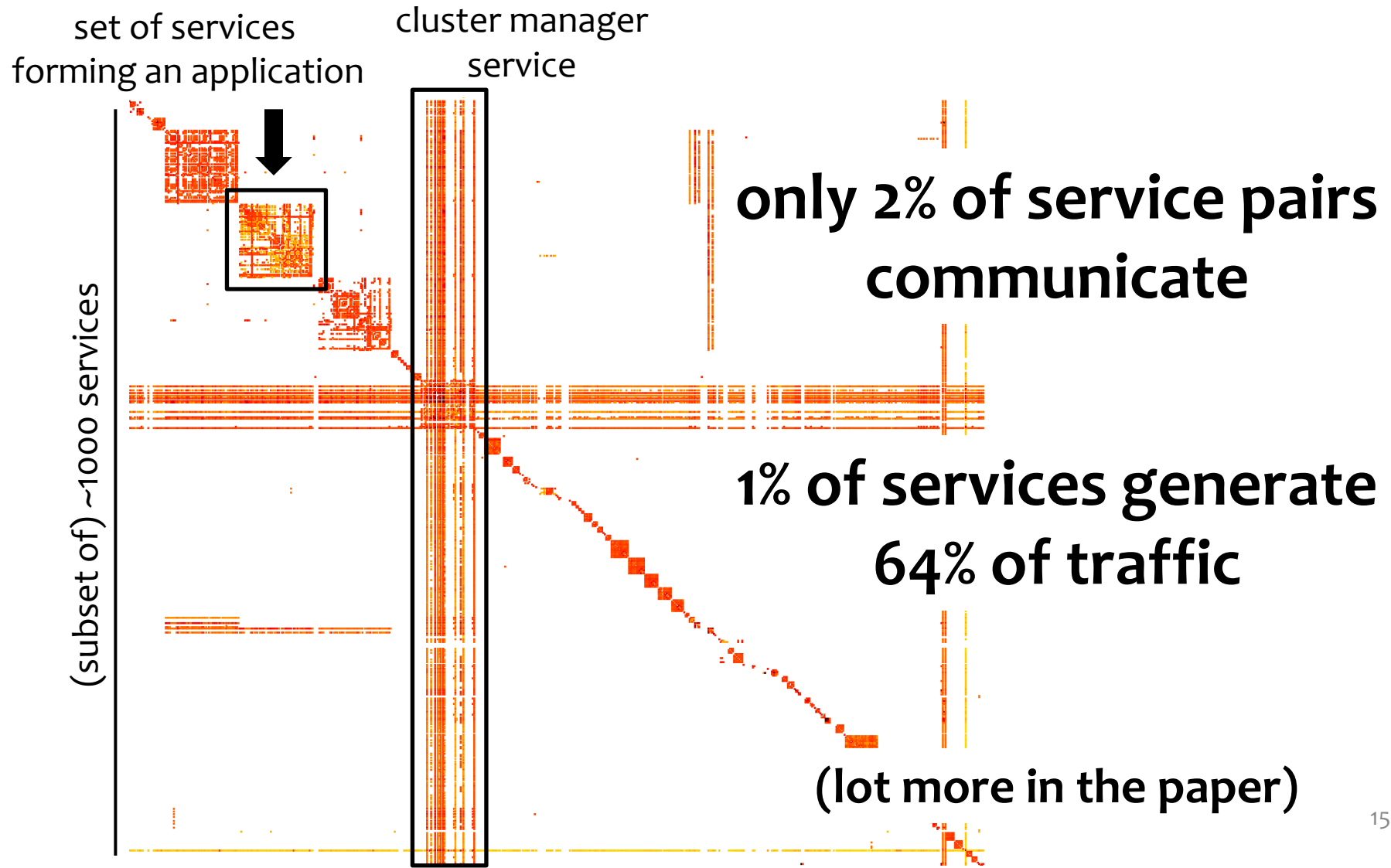
Optimization framework

- FT + #M

- FT + BW + #M

Evaluation

# Service communication matrix is very sparse and skewed



# Outline

Why is it difficult?

Traffic analysis

Optimization framework

- FT + #M

- FT + BW + #M

Evaluation





# optimizing FT and #M

Spread machines across all fault domains

- FTC negatively correlated to worst-case survival

Convex optimization

$$\text{FTC} = \sum_s c_s \sum_f w_f \cdot z_{s,f}^2$$

service weight
fault domain weight

number of machines of service  $s$  in domain  $f$

Advantages of convex cost function

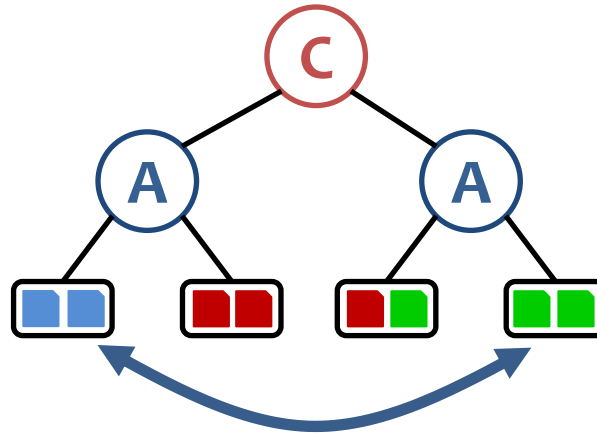
- local actions lead to improvement of global metric
- directly considers #M

**FT**

# machine swap as a basic move

Keeps the current allocation feasible

- doesn't change number of machines per service



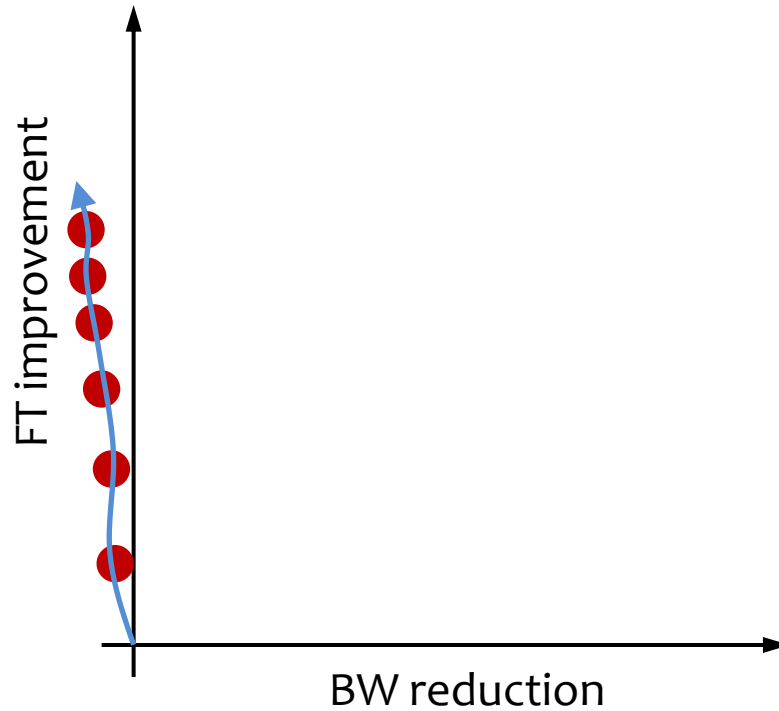
Steepest descent swap = largest reduction in cost

Only evaluate a small, random set of swaps

- symmetry => many “good” swaps exist



# path of steepest descent



# FT+BW

# Optimizing FT, BW, and #M

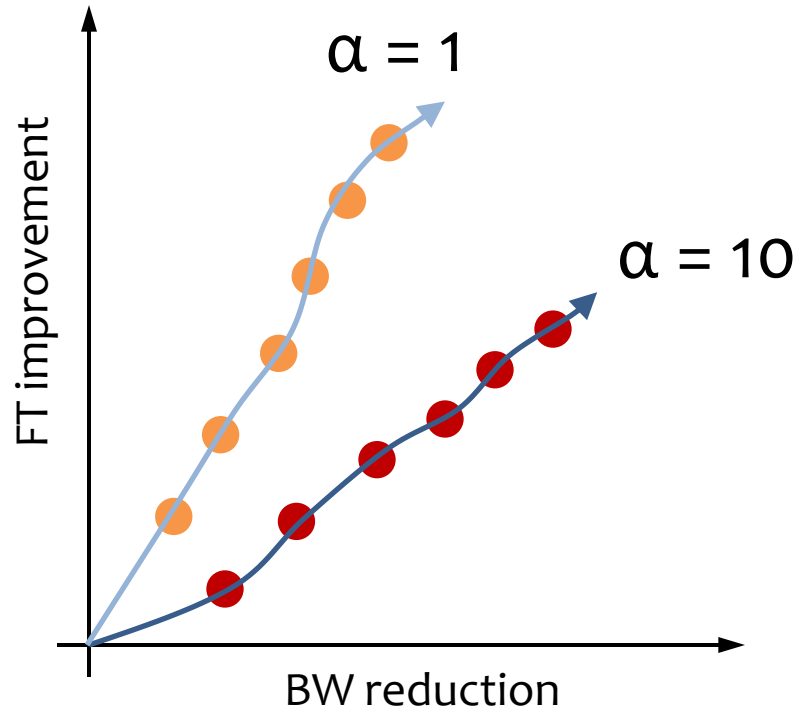
Steepest descent on  $FTC + \alpha BW$

- non-convex
- no guarantees on reaching optimum

$\alpha$  determines the FT-BW trade-off

**FT+BW**

# path of steepest descent



# Benchmark algorithm

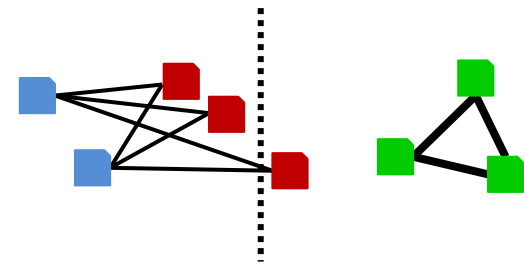
**cut**

**FT+BW**

k-way minimum graph cut

- optimizes BW only
- ignores #M

machine communication graph



k-way min cut

followed by steepest descent on FT+BW

# Outline

Why is it difficult?

Traffic analysis

Optimization framework

- FT + #M

- FT + BW + #M

Evaluation

# Evaluation setup

## Simulations based on 4 production clusters

- services + machine counts
- network topology
- fault domains
- network trace from pre-production cluster

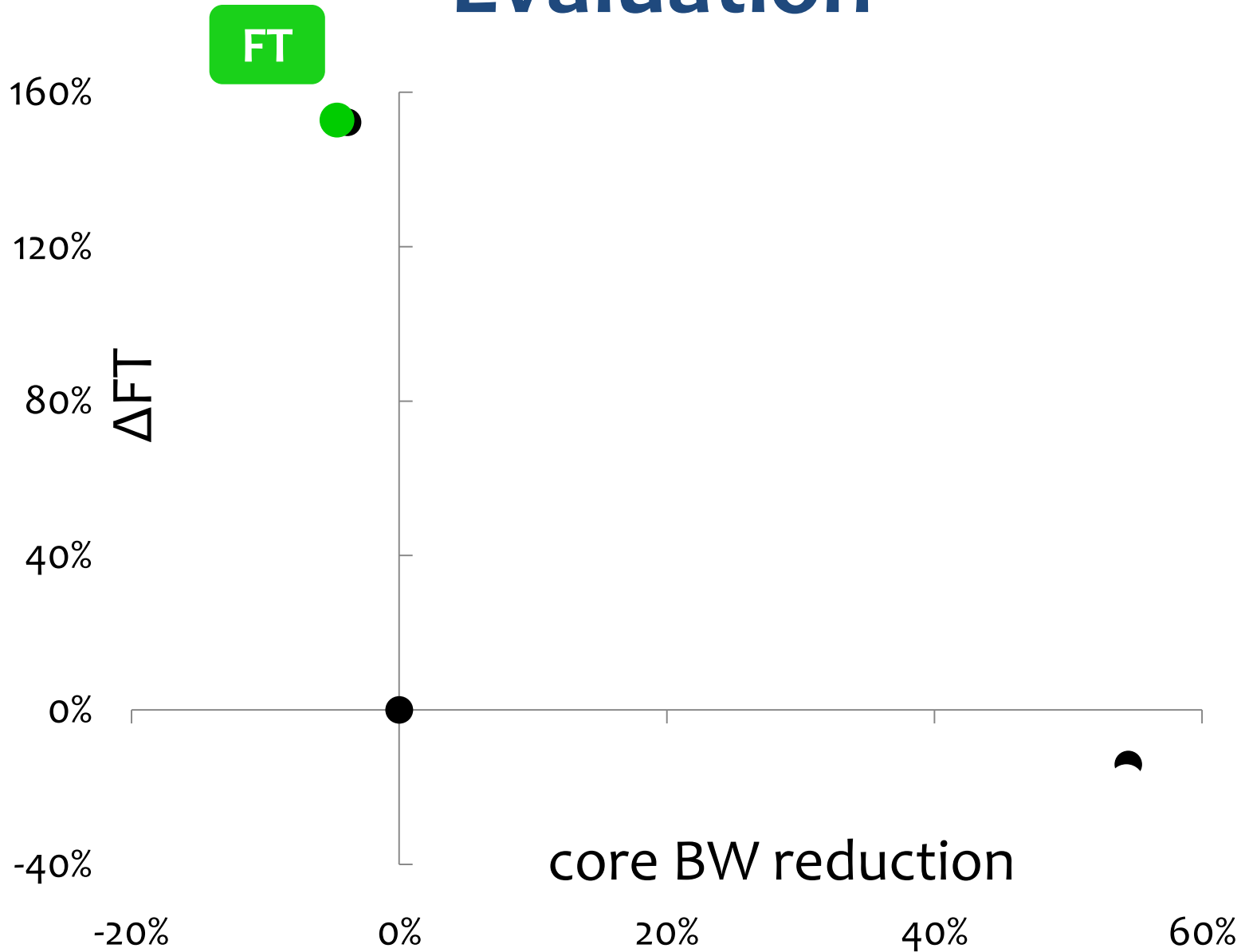
## Metrics relative to initial allocation

- don't know actual optimum

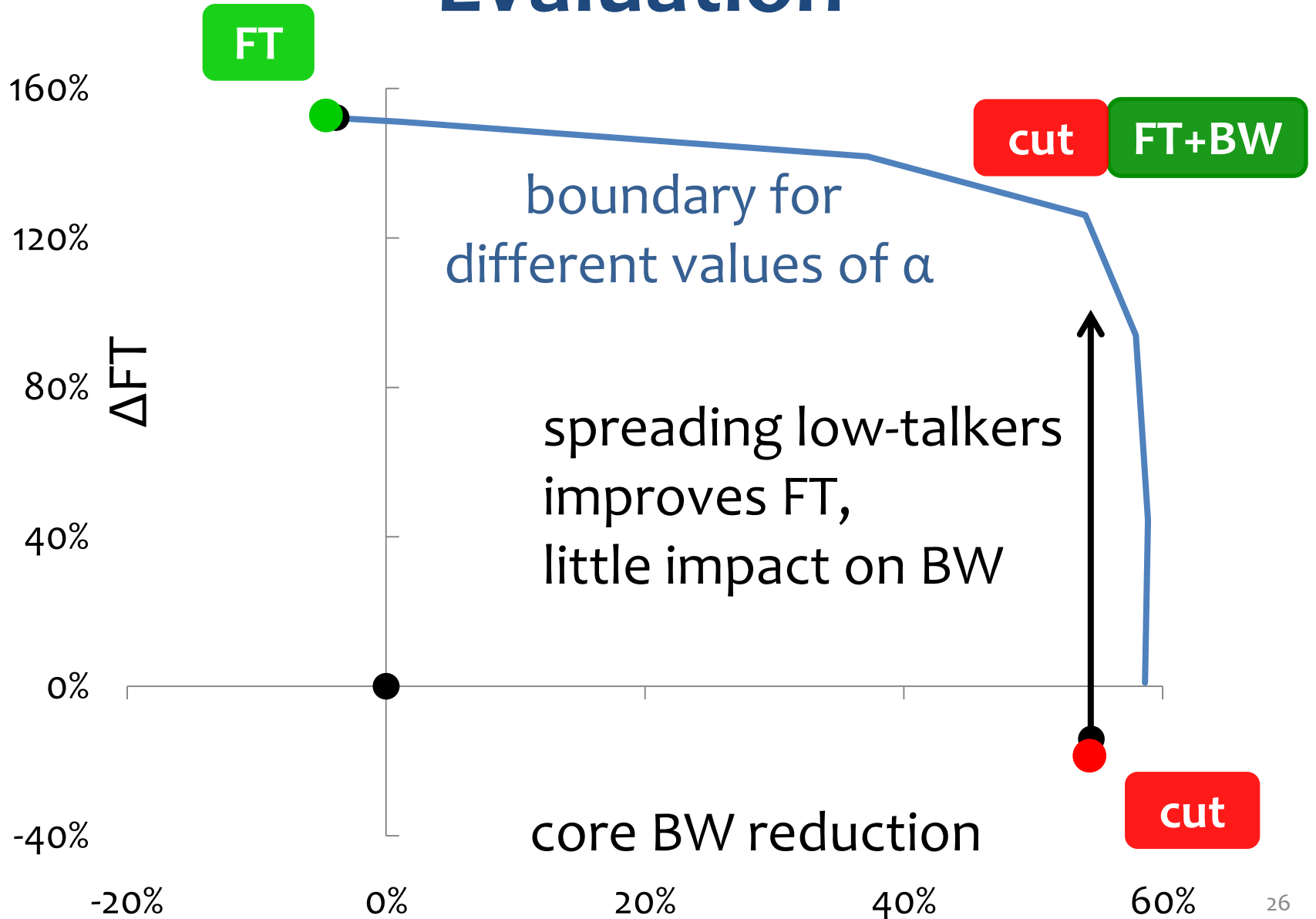
Choosing next swap takes seconds to a minute



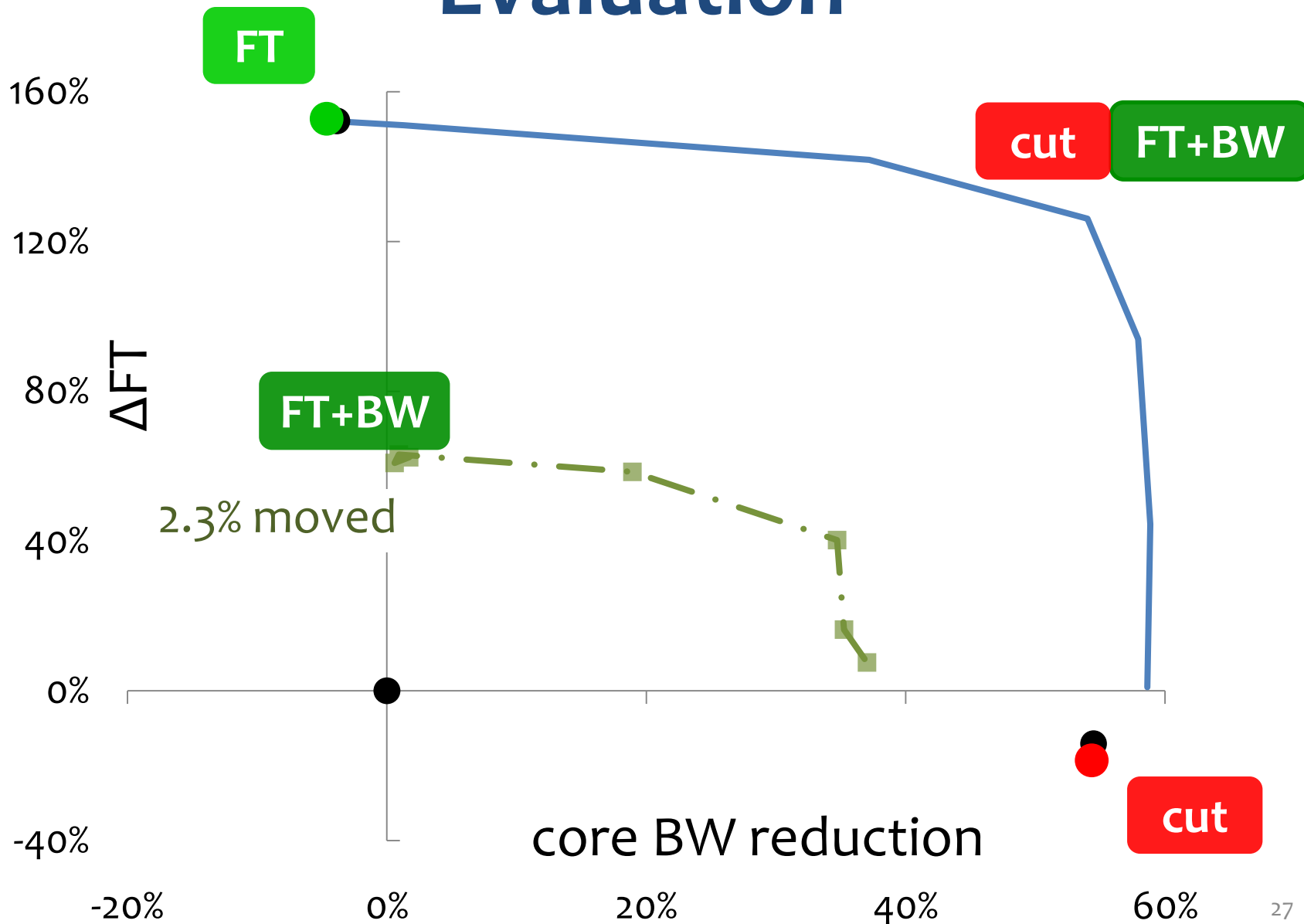
# Evaluation



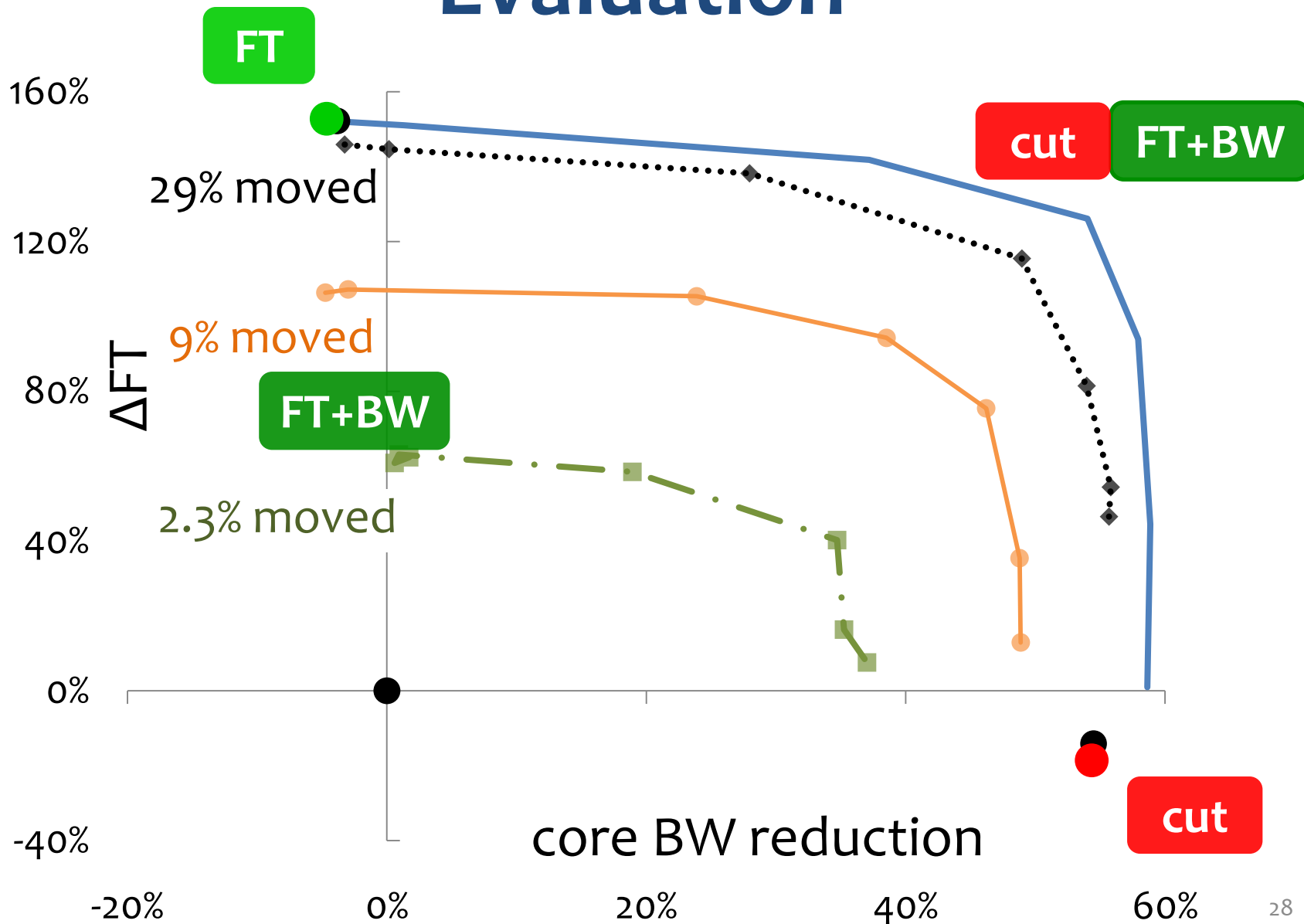
# Evaluation



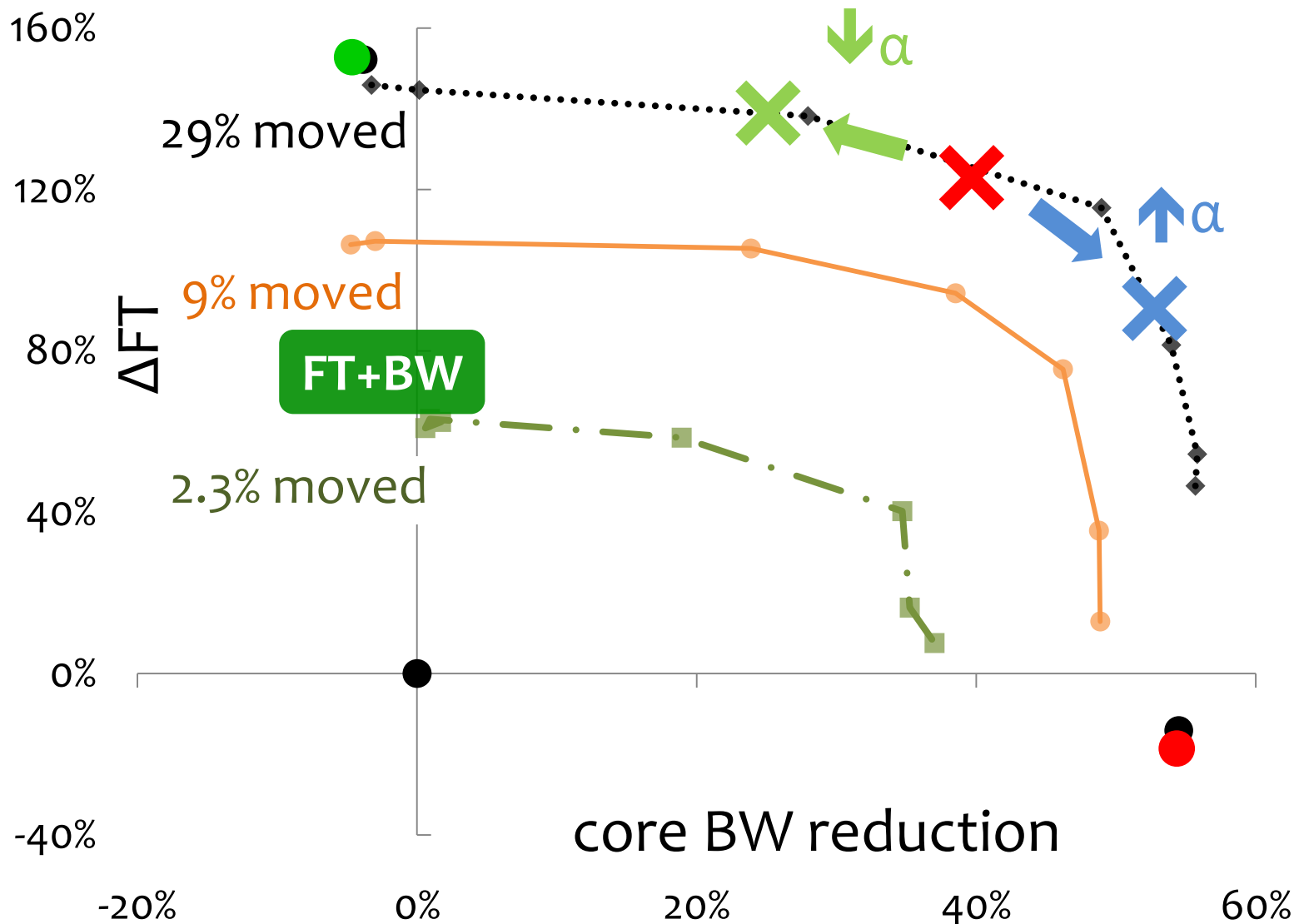
# Evaluation



# Evaluation



# $\alpha$ changes the FT-BW tradeoff



# Summary

Trade-off between fault tolerance and bandwidth

- algorithm that achieves improvement in both

Improvements (across 4 production datacenters)

- FT: 40% – 120%
- BW: 20% – 50%
- partially deployed in Bing

Key insights

- approximate NP-hard problem using convex optimization
- lot of symmetry in search space
- sparse and skewed communication matrix







# Extensions

Hard constraints on FT, BW, #M

- e.g., pick a few services with  $FT > 80\%$

Hierarchical BW optimization on agg switches

Applies to fat-tree networks

# Main observations

Most traffic generated by few services (pairs)

➔ spread low-talkers to improve fault-tolerance

Complex, overlapping fault domains

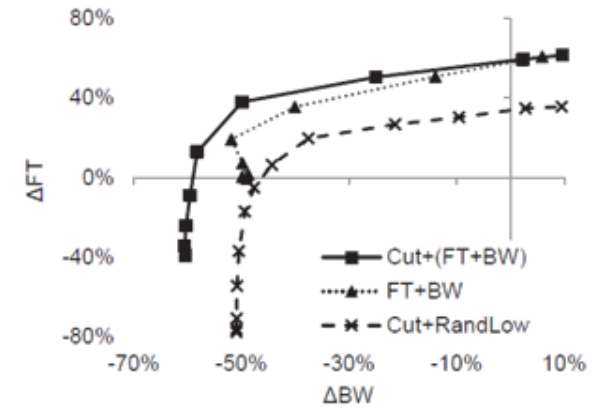
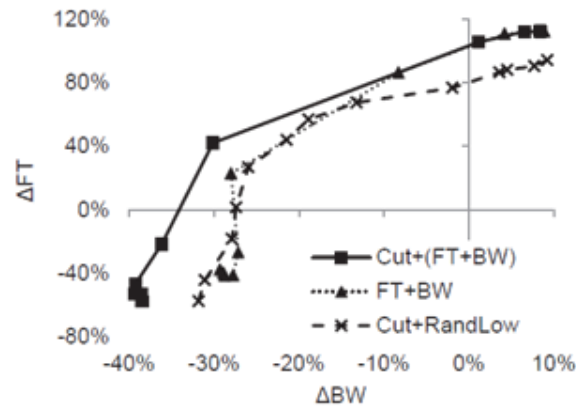
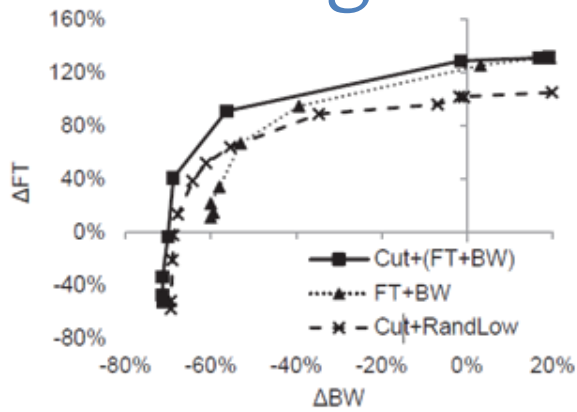
– hierarchical network fault-domains

– power fault domains not aligned with network

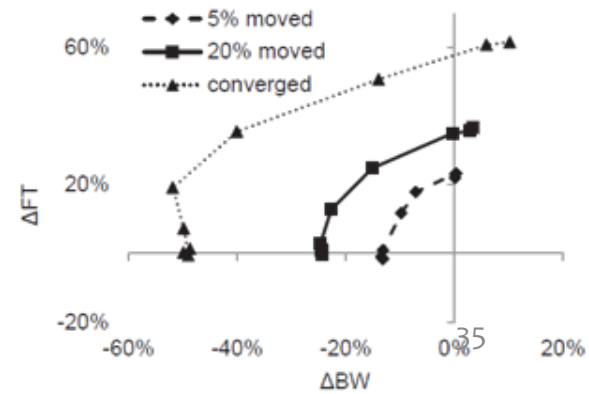
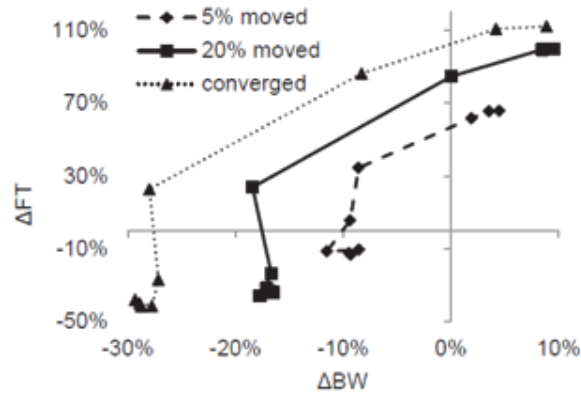
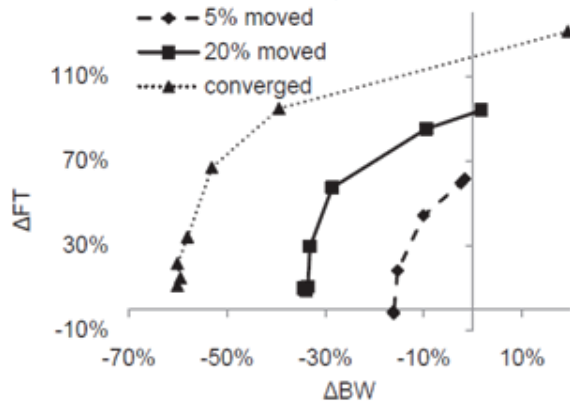
➔ **cell:** set of machines with identical fault domains

# Evaluation

## Moving most of machines



## Moving only fraction of machines



# Our optimization framework

## Cost function considers FT and BW

- both problems NP-hard and hard to approximate
- non-convex

## Cut + FT + BW:

1. minimum k-way cut of communication graph
  - reshuffles all machines
2. gradient descent moves using machine swaps

## FT + BW:

1. only machine swaps
  - only moves small fraction of machines

# Conclusion

## Study of communication patterns of Bing.com

- sparse communication matrix
- very skewed communication pattern

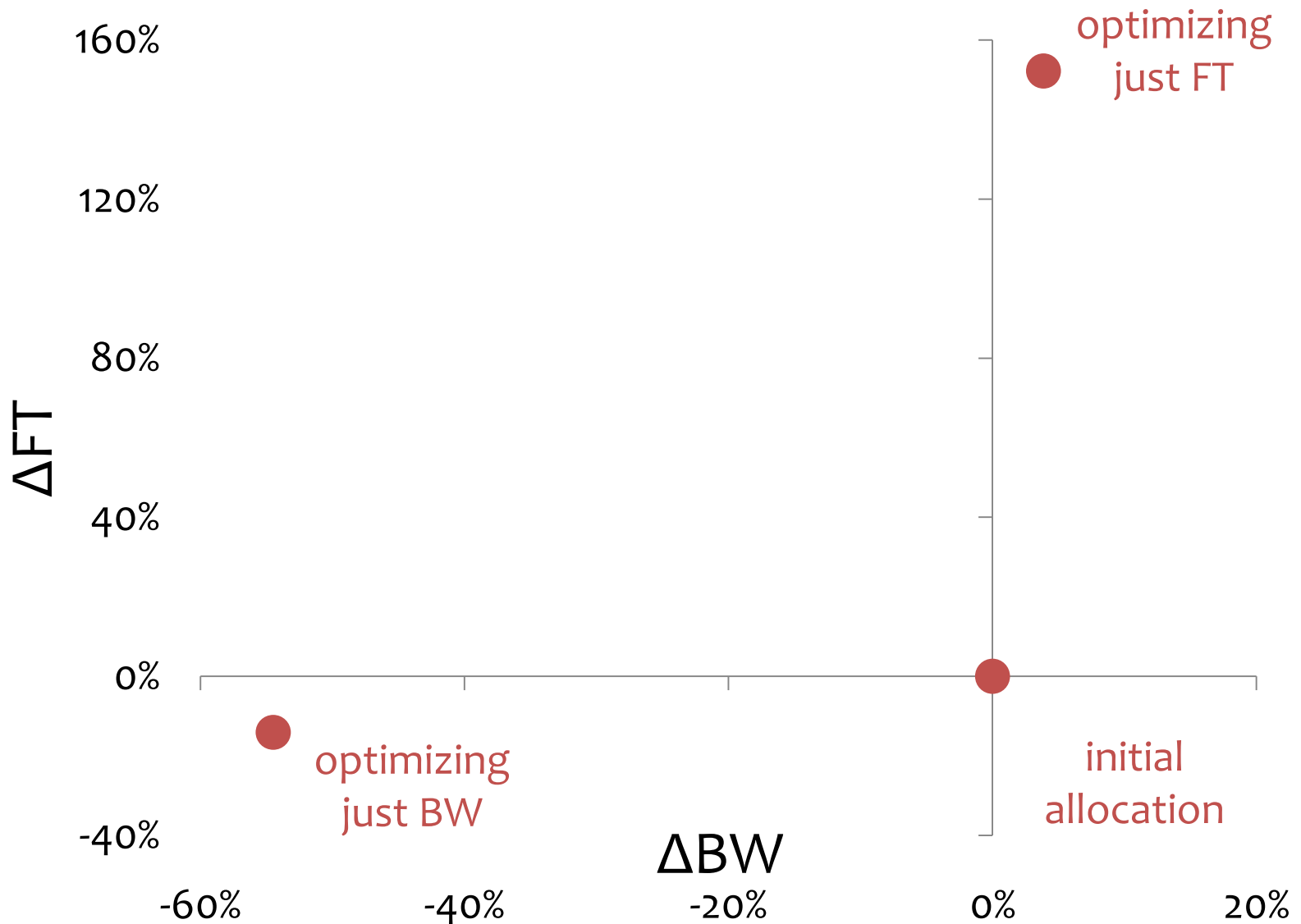
## Principled optimization of both BW and FT

- exploits communication patterns
- can handle arbitrary fault domains

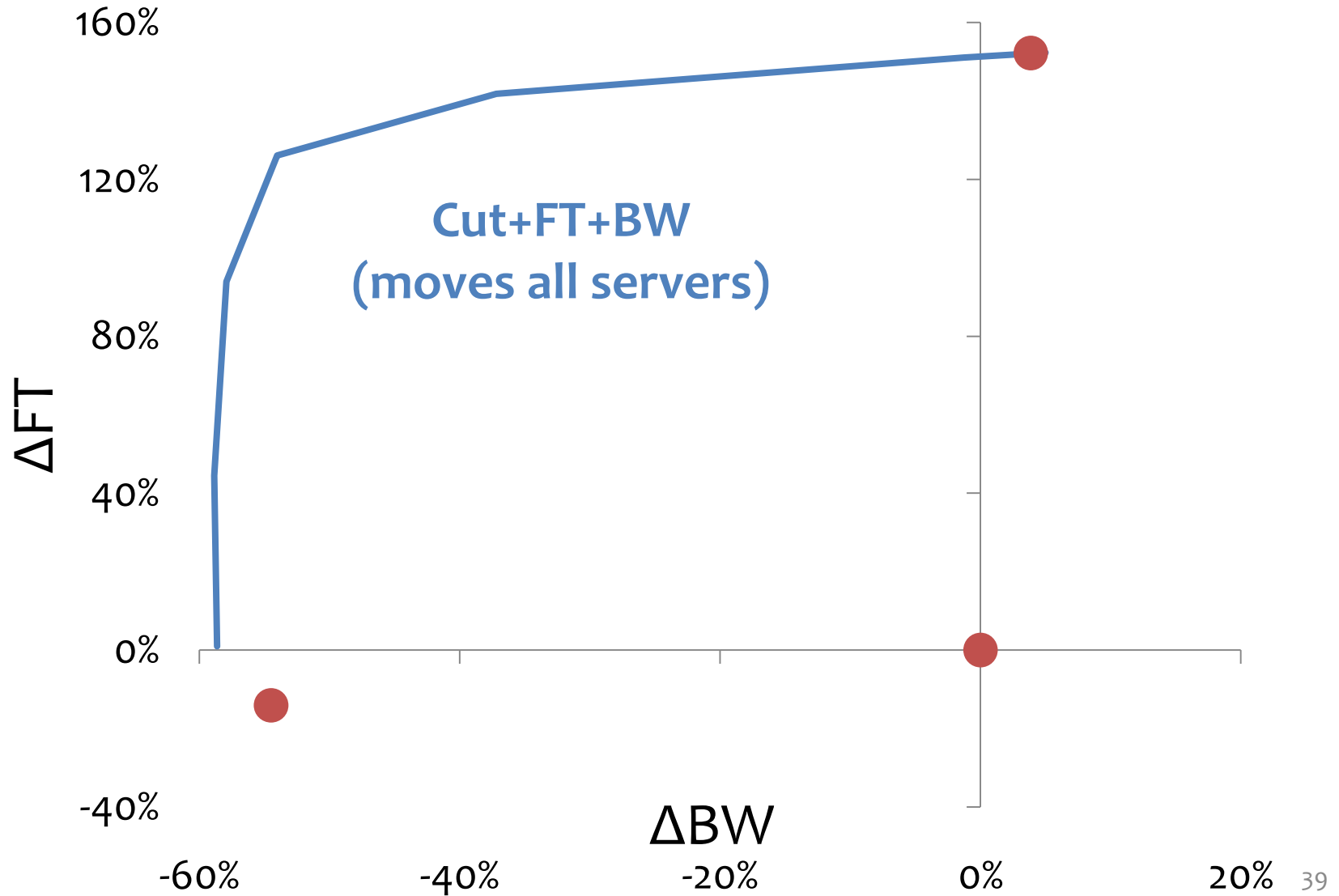
Reduction in BW: 20 – 50%

Improvement in FT: 40 – 120%

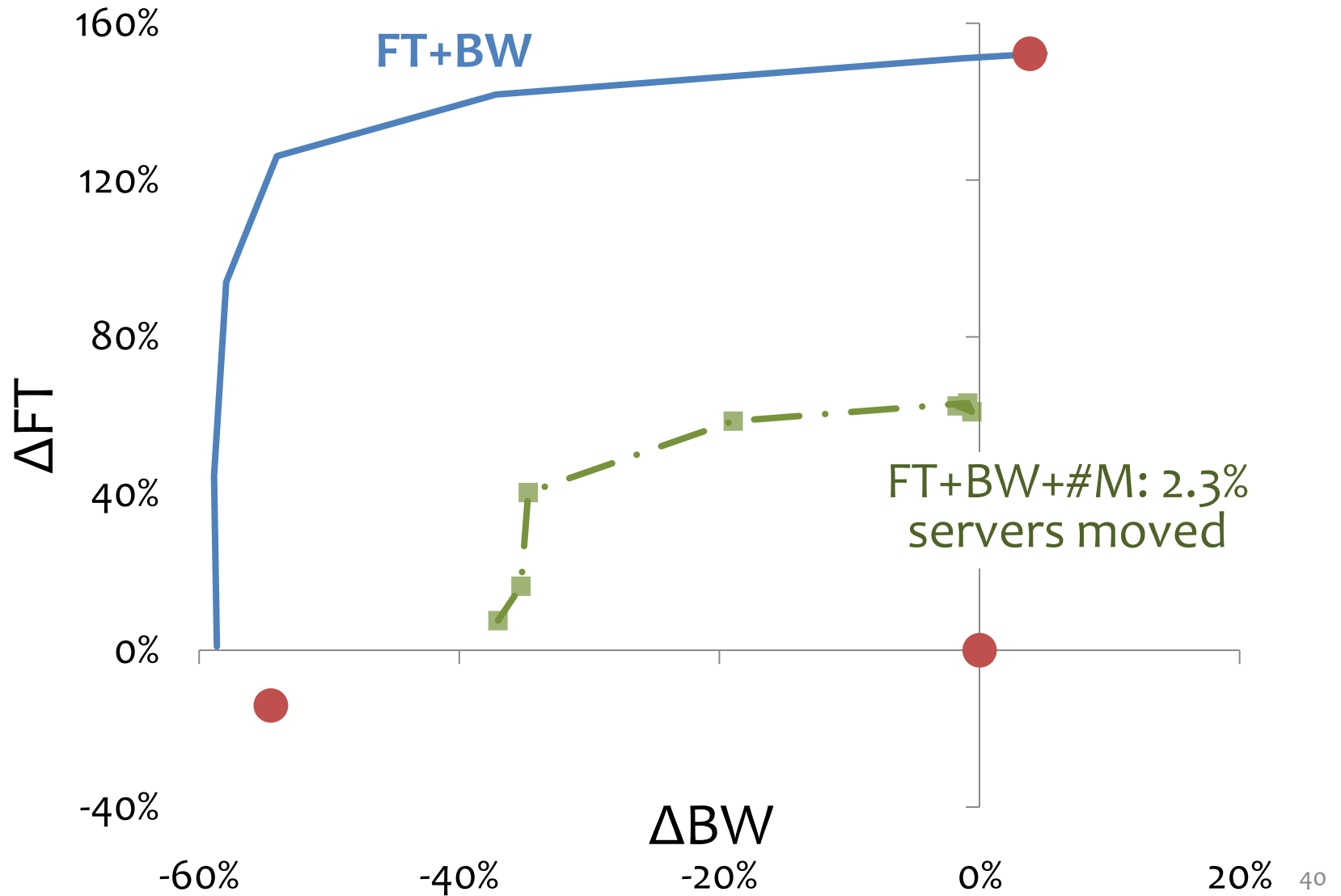
# Evaluation (1 datacenter)



# Evaluation

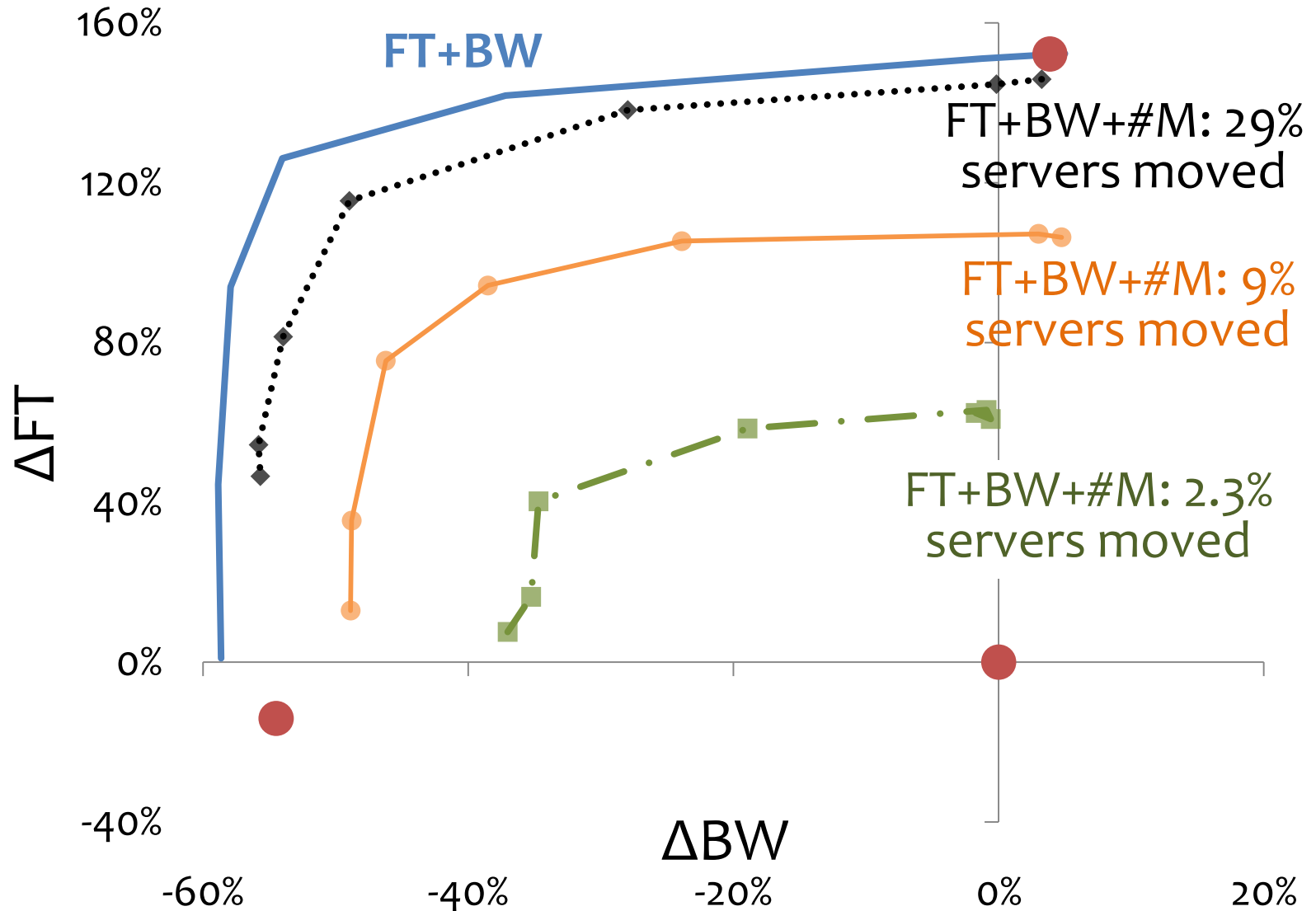


# Evaluation

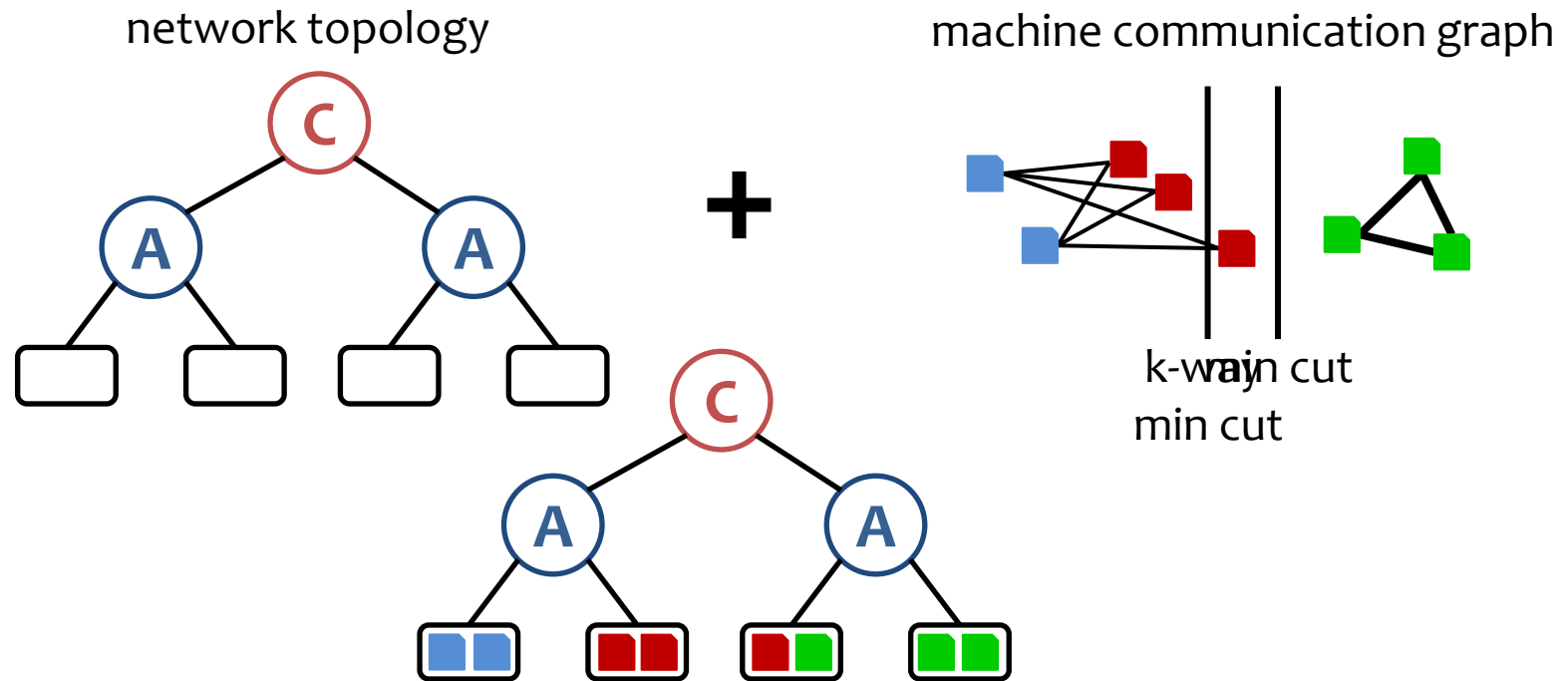




# Evaluation



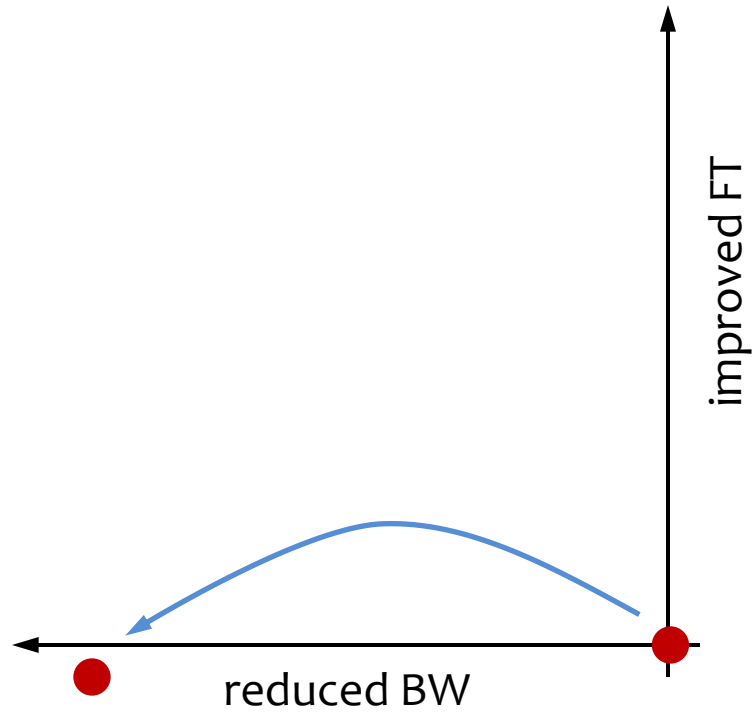
# k-way graph cut



## k-way min graph cut

- ignores #M: reshuffles almost all machines
- ignores FT: can't be easily extended

# BW k-way graph cut



# Scaling algorithms to large datacenters

Only evaluate a small, random set of swaps

- symmetry => many “good” swaps exist

Cell = set of machines with same fault domains

Reduce size of communication graph for cut



# **FT** **BW** cut + steepest descent

Step 1: min-cut

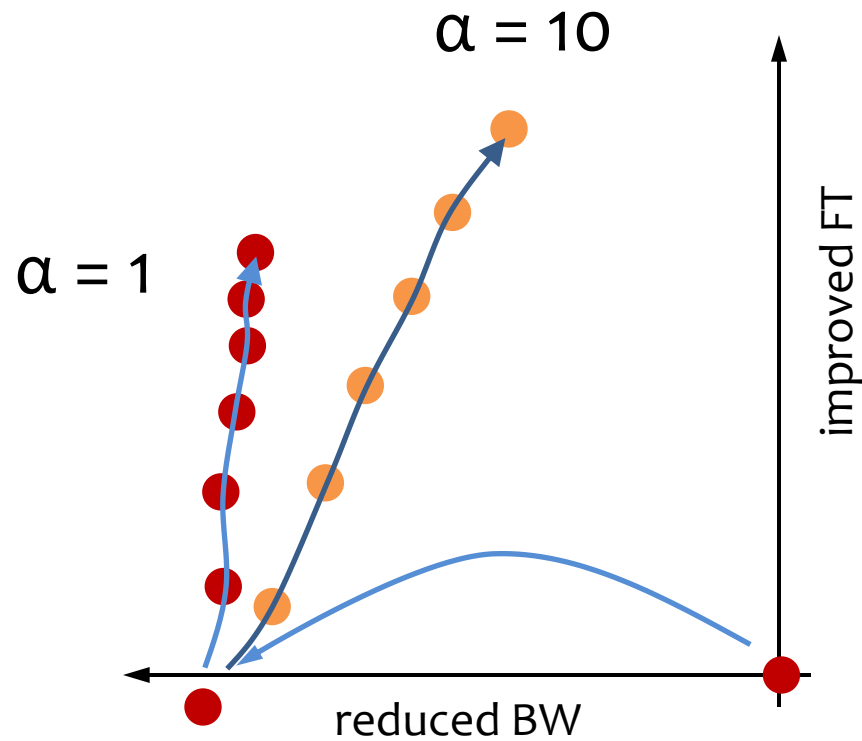
- optimizes BW

Step 2: steepest descent on  $FTC + \alpha BW$

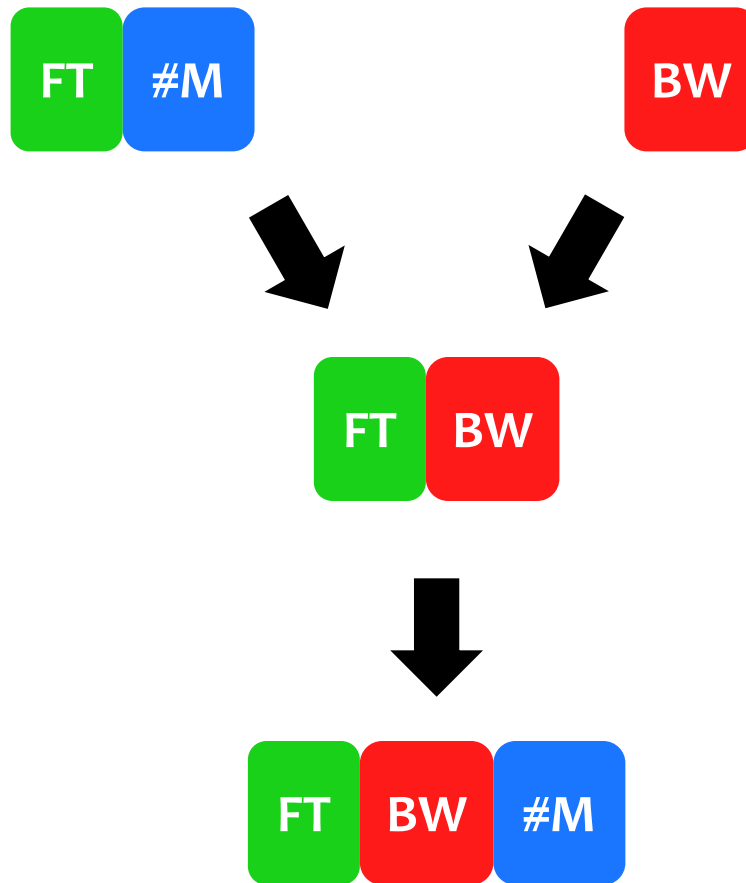
- non-convex
- no guarantees on reaching optimum
- $\alpha$  determines the trade-off

Reshuffles all machines

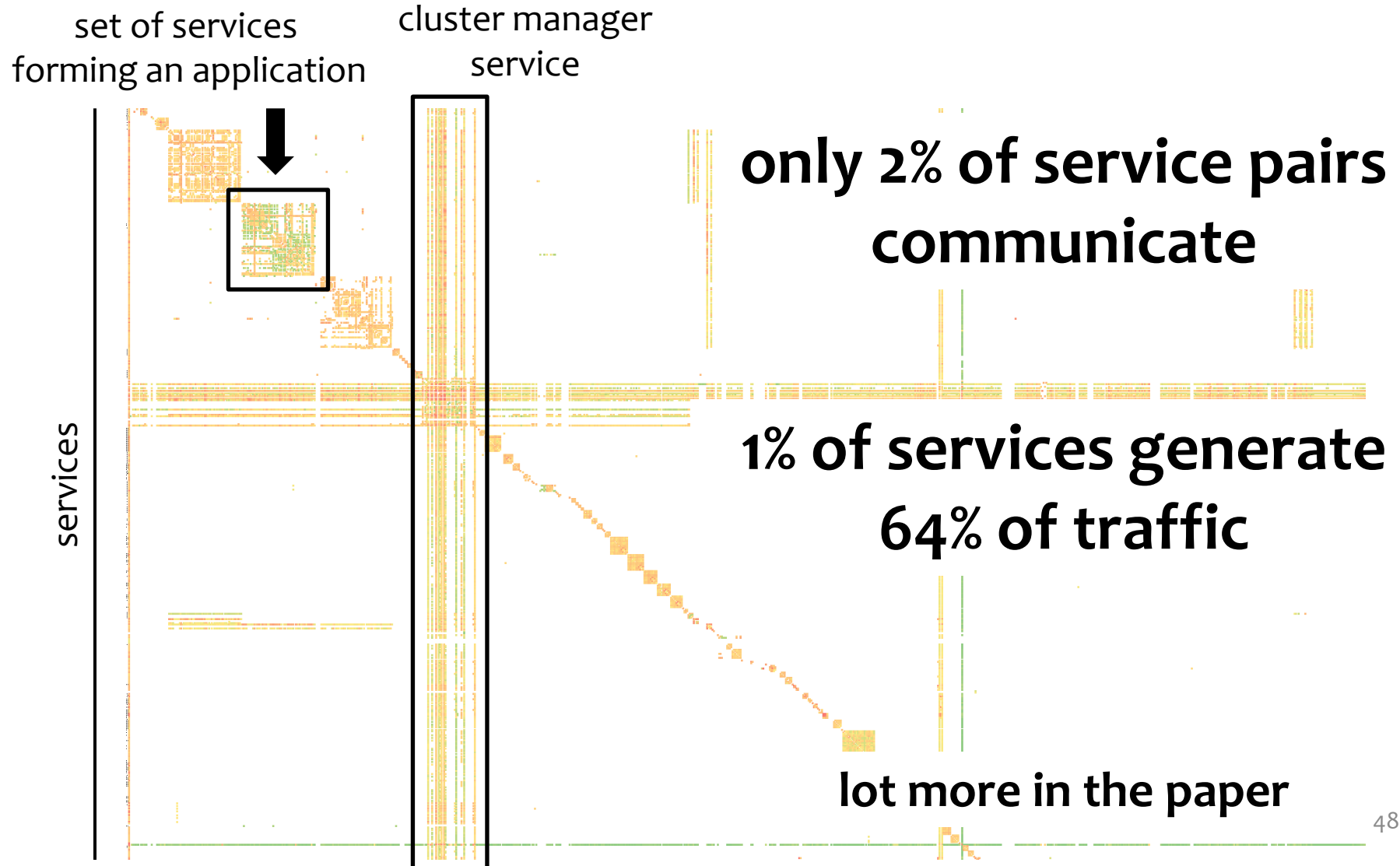
# FT BW cut + steepest descent



# Properties of allocation algorithms



# Service communication matrix is very sparse and skewed





# Which metrics matter?

## FT: fault tolerance

- service should survive infrastructure failures
- failures despite redundancy

## BW: bandwidth usage

- reduce usage on constrained links
- lower cost of infrastructure

## #M: number of moves

- moving some servers is expensive
- want incremental allocation