# CONTENT-AWARE LOAD BALANCING FOR DISTRIBUTED BACKUP

Fred Douglis[1], Deepti Bhardwaj[1], Hangwei Qian[2], and Philip Shilane[1]

[1]EMC          [2]Case Western Reserve University

# Starting Point

- **Deduplicating** disk-based backup storage
  - Variable, content-defined chunks
    - Strong hashes of content to find duplicates

**EMC²**

# Starting Point

- **Deduplicating** disk-based backup storage
  - Variable, content-defined chunks
    - Strong hashes of content to find duplicates
  - Focused on making full backups after the first one use minimal extra disk space
    - Internal deduplication – duplicates from multiple copies of the same file from the same source
    - Unchanged files dedupe trivially, while chunk-level deduplication catches changes scattered within large regions of unchanged content

**EMC²**

# Starting Point

- **Deduplicating** disk-based backup storage
  - Variable, content-defined chunks
    - Strong hashes of content to find duplicates
  - Focused on making full backups after the first one use minimal extra disk space
    - Internal deduplication – duplicates from multiple copies of the same file from the same source
    - Unchanged files dedupe trivially, while chunk-level deduplication catches changes scattered within large regions of unchanged content
  - Deduplication can avoid sending the data at all
    - Send the hashes and then only send new chunks

**EMC²**

# Starting Point

- ## Deduplicating disk-based backup storage
  - Variable, content-defined chunks
    - Strong hashes of content to find duplicates
  - Focused on making full backups after the first one use minimal extra disk space
    - Internal deduplication – duplicates from multiple copies of the same file from the same source
    - Unchanged files dedupe trivially, while chunk-level deduplication catches changes scattered within large regions of unchanged content
  - Deduplication can avoid sending the data at all
    - Send the hashes and then only send new chunks
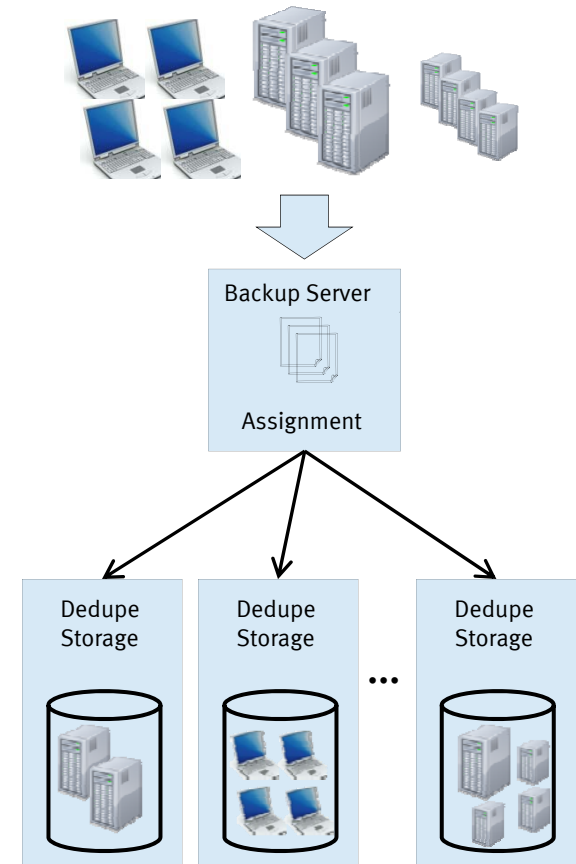  - Technology now common in backup products

**EMC²**

# Problem Statement

- Large-scale IT environment
  - Hundreds or thousands of systems ("clients") to backup
  - Many backup appliances to send the data

**EMC²**

# Problem Statement

- Large-scale IT environment
  - Hundreds or thousands of systems ("clients") to backup
  - Many backup appliances to send the data

- Impact of deduplication
  - Affinity: send the same client to the same appliance so it will deduplicate well
    - Moving it to a new system will cause everything to be written again
  - Overlap: benefit from sending *similar* systems to the same backup appliance
    - "External" deduplication, spanning clients

**EMC²**

# Problem Statement

- Large-scale IT environment
  - Hundreds or thousands of systems ("clients") to backup
  - Many backup appliances to send the data

- Impact of deduplication
  - Affinity: send the same client to the same appliance so it will deduplicate well
    - Moving it to a new system will cause everything to be written again
  - Overlap: benefit from sending *similar* systems to the same backup appliance
    - "External" deduplication, spanning clients

Backup Server

Assignment

Dedupe Storage

Dedupe Storage

...

Dedupe Storage

Simple approach: cluster clients by type

EMC²

# Benefits of Overlap

- Co-locating duplicate content
  - Reduces capacity requirements
    - May take a host from being overloaded to highly loaded, or highly loaded to moderately
  - Reduces throughput requirements
    - Duplicate copies in later clients' **first full** are skipped
    - Ongoing transfers benefit only if identical content being written to multiple hosts during a backup interval

- <span style="color:red">Deduplication changes traditional backup administration</span>
  - Backup devices are not all created equal
    - They're not all identical <u>tapes</u>
  - There is a "stickiness" to the assignment in order to benefit from savings
  - But sometimes data migration benefits outweigh costs

**EMC²**

# Benefits of Overlap

- Co-locating duplicate content
  - Reduces capacity requirements
    - May take a host from being overloaded to highly loaded, or highly loaded to moderately
  - Reduces throughput requirements
    - Duplicate copies in later clients' **first full** are skipped
    - Ongoing transfers benefit only if identical content being written to multiple hosts during a backup interval

- Deduplication changes traditional backup administration
  - Backup devices are not all created equal
    - They're not all identical <u>tapes</u>
  - There is a "stickiness" to the assignment in order to benefit from savings
  - But sometimes data migration benefits outweigh costs

## Where do we put clients and when do we have to give in and move them?

EMC²

# Goals

- Capacity allocation
  - Send data to backup appliances in the best way to fit them within constraints
  - Balanced load
  - Content-aware for best deduplication

**EMC²**

# Goals

- **Capacity** allocation
  - Send data to backup appliances in the best way to fit them within constraints
  - Balanced load
  - Content-aware for best deduplication

- **Performance** (throughput)
  - Support many backup streams simultaneously
    - Avoid overloading any individual appliances
  - Increased deduplication reduces overhead on network and appliance

**EMC²**

# Use Cases

- ## Sizing and deployment
  - Figure out requirements (and assignments) from "clean slate"

- ## First assignment
  - Given a set of clients and appliances, determine best assignments

- ## Reconfiguration
  - Adjust when clients or appliances are added or removed, or load shifts

- ## Disaster recovery & replication
  - Select mappings of appliances onto other appliances for off-site replication

**EMC²**

# Approach

- Minimize a <span style="color:red">utility function</span>
  - "Cost" of a configuration is a function of capacity utilization and performance requirements
    - Compare costs directly to identify best configuration
  - Lots of tradeoffs
    - E.g., migrate a client to a new appliance to reduce capacity overload, but pay a penalty for data movement

- Identify overlap
  - Sample fingerprints for each client
  - Find cases of "significant" overlap
    - Ignore the rest

**EMC²**

# How Much Overlap is There?

- Many systems will have little or no overlap

- Some systems will have similar overlap with many other systems, so picking one in particular has no advantage

- Want to identify special affinity in cases of high overlap among 2, or few, hosts

- Studied 21 hosts from saved workstation backups and live systems
  - One host with 50% overlap with another and almost 25% **additional** overlap with a third
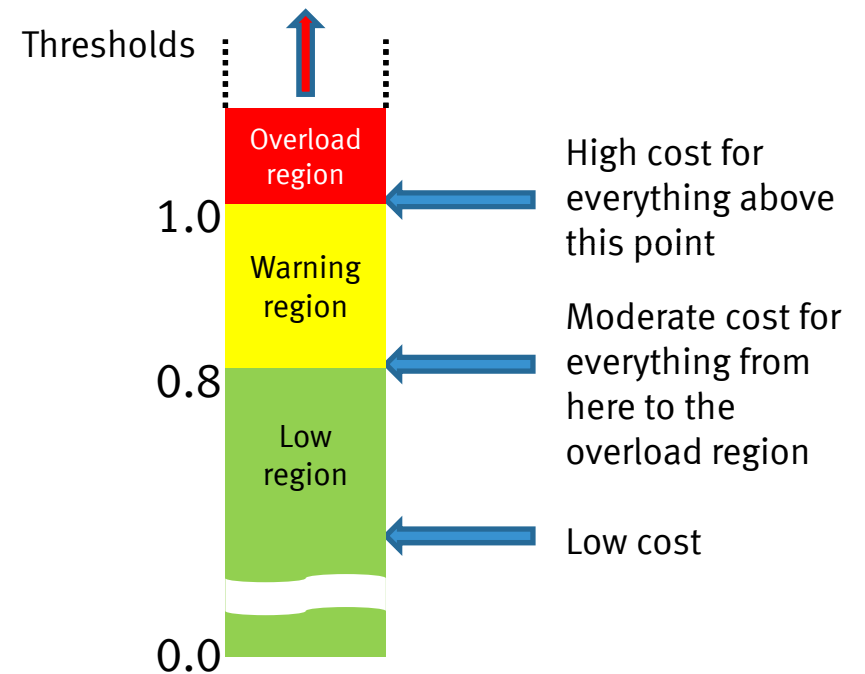
- Virtual machine images particularly likely to have high overlap

**Deduplication (%)**

host16

host21

host20

host16

host20

host4

■ Available deduplication

■ 2nd match

■ Best Match

■ Widely deduped

host1 host2 host3 host4 host5 host6 host7 host8 host9 host10 host11 host12 host13 host14 host15 host16 host17 host18 host19 host20 host21

**EMC²**

# Cost Calculation

- In total, the cost for a given configuration is the sum of:
  - A small, weighted penalty for imbalance in capacity or throughput

# Cost Calculation

- In total, the cost for a given configuration is the sum of:
    - A small, weighted penalty for imbalance in capacity or throughput
    - A stepped penalty for exceeding thresholds in capacity or throughput

Thresholds

Overload region

1.0

Warning region

0.8

Low region

0.0

High cost for everything above this point

Moderate cost for everything from here to the overload region

Low cost

**EMC²**

# Cost Calculation

- In total, the cost for a given configuration is the sum of:
  - A small, weighted penalty for imbalance in capacity or throughput
  - A stepped penalty for exceeding thresholds in capacity or throughput
  - A small penalty for migrating off an existing appliance

# Cost Calculation

- In total, the cost for a given configuration is the sum of:
  - A small, weighted penalty for imbalance in capacity or throughput
  - A stepped penalty for exceeding thresholds in capacity or throughput
  - A small penalty for migrating off an existing appliance
  - A <span style="color:red">very large penalty</span> for each client that does not "<span style="color:red">fit</span>" on its appliance
    - *In our experiments presented today, this penalty is the dominant cost. Above 1000 means "overload" and below it means "fit"*
    - *Smaller penalties are used to pick among plausible choices*

- (A more formal definition appears in the paper)

**EMC²**

# Algorithms

- Compare "intelligent" assignment to brute force such as round-robin or random
  - All the brute force approaches quite fast

- Random
  - Pick arbitrary assignments. If random selection is full, iterate to find new appliance.
  - Compute cost of configuration
  - Repeat N times and take best result

- Round-robin
  - Assign to each appliance in turn
  - Skip a "full" appliance to find one with available capacity if possible

- Bin packing
  - Assign based on size from largest to smallest (less likely to overflow)

- Simulated annealing
  - Shuffles assignments from the current "best position" to try and improve the cost

- The first three take any existing assignments as a given; only annealing will migrate a client

- Generally, all work well under low load; annealing can adapt better to overload

EMC²

# Annealing Example

utilization



swap

move

**EMC²**

# Evaluation (Simulations)

# Incremental Assignment Experiment

- Define a number of clients of fixed size: small, medium, large, 20 per iteration

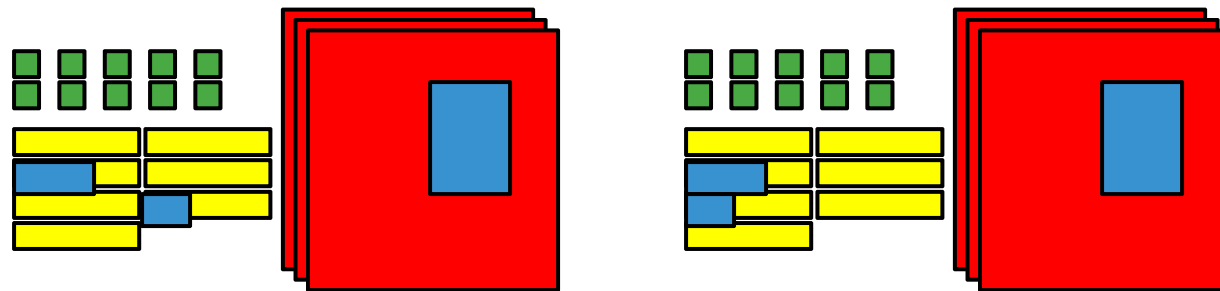| 2 | 0 | G | B | |

100 GB

2 TB

**EMC²**

# Incremental Assignment Experiment

- Define a number of clients of fixed size: small, medium, large, 20 per iteration

- Repeatedly put a set of clients into system and assign to appliances
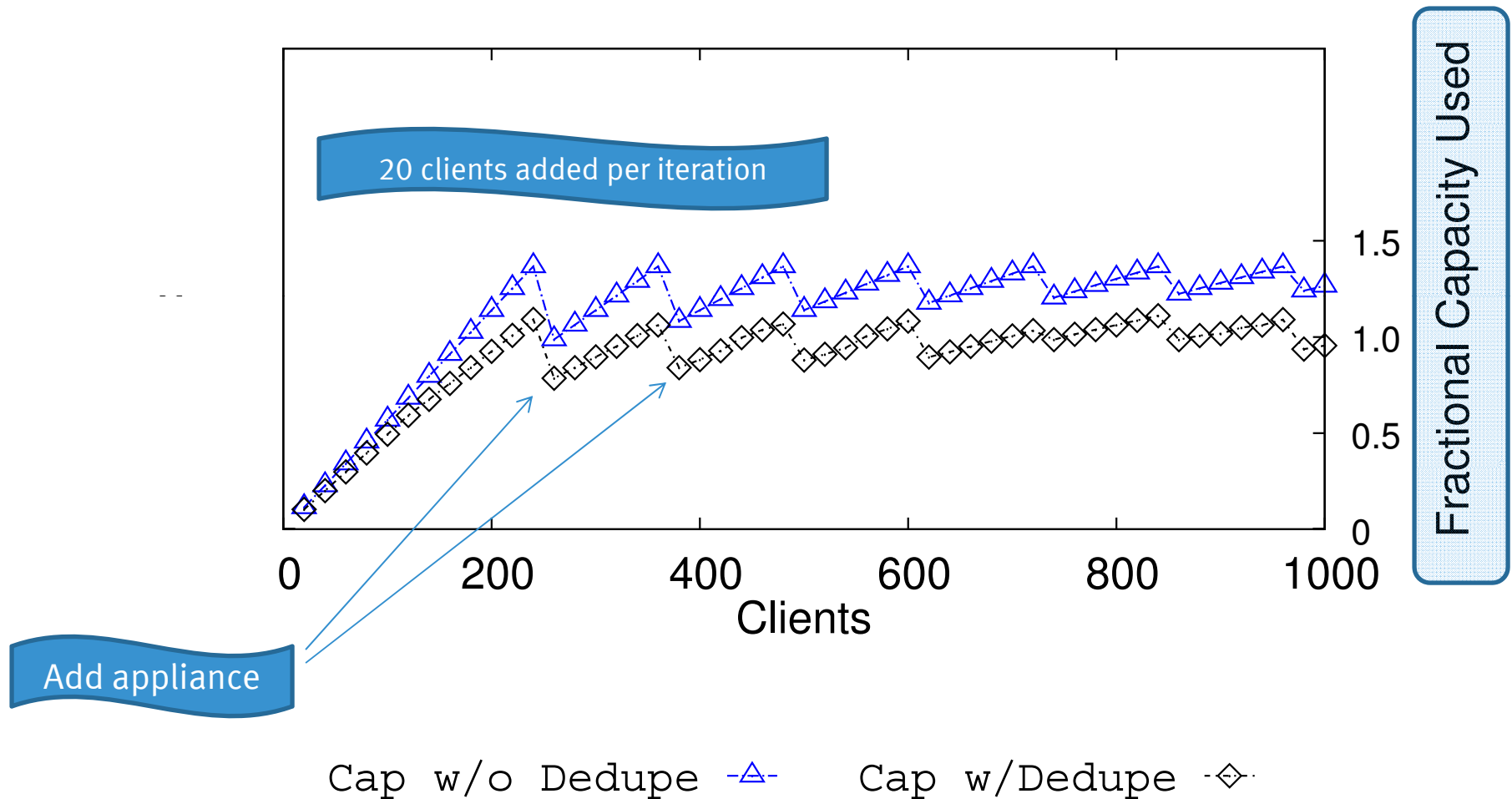  - Better dedupe within a class than across

# Incremental Assignment Experiment



- **Periodically add a new appliance** to increase capacity
  - At the same time, forget 1/3 of existing assignments (so some assignments have a penalty for movement and some don't)
  - Especially high dedupe with the corresponding client from other iterations – stress overlap affinity

- If new load outpaces capacity, high cost.  If the new appliance is added to keep up with added load, low cost.
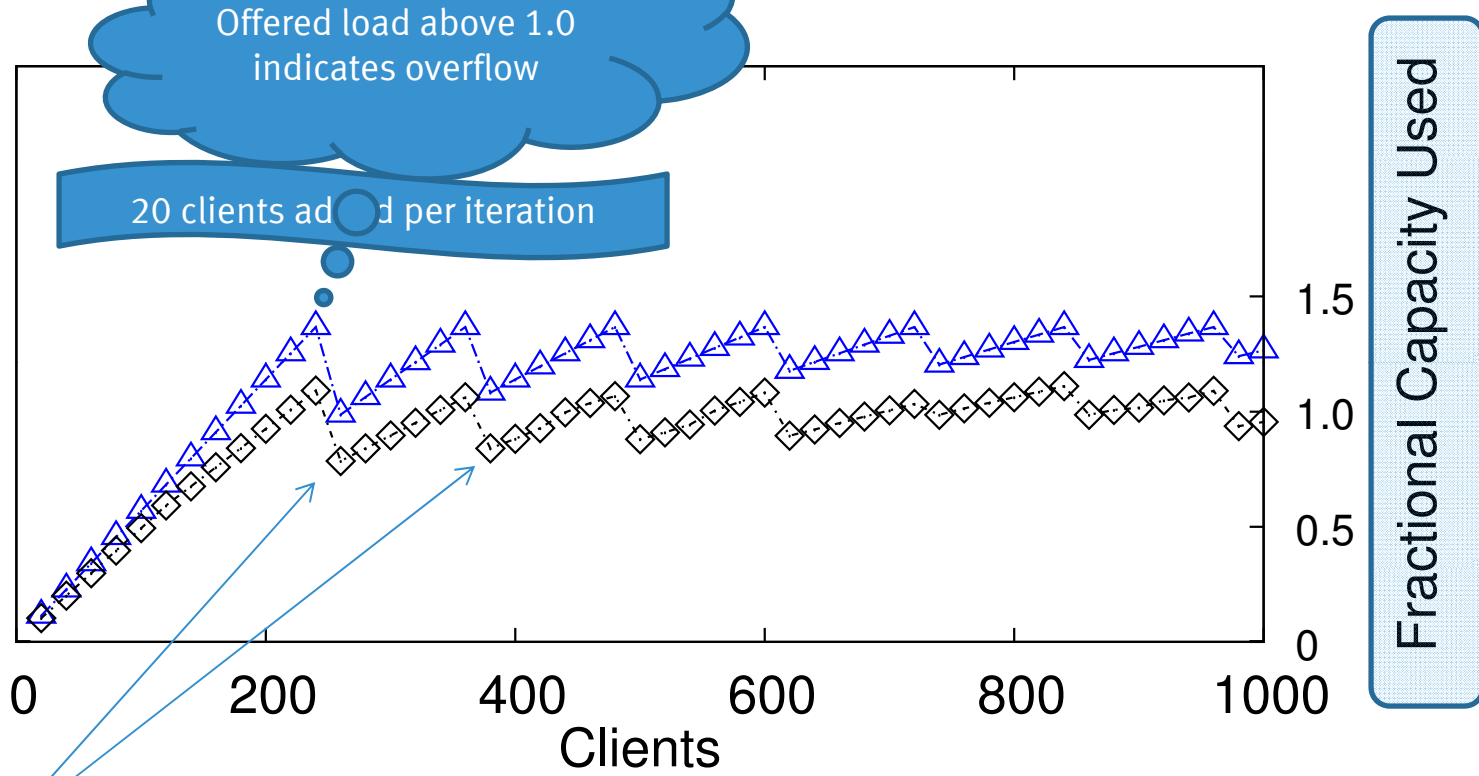
**EMC²**

# Incremental Assignment Experiment

- Define a number of clients of fixed size: small, medium, large, 20 per iteration

- Repeatedly put a set of clients into system and assign to appliances
  - Better dedupe within a class than across

- <span style="color:red">Periodically add a new appliance</span> to increase capacity
  - At the same time, forget 1/3 of existing assignments (so some assignments have a penalty for movement and some don't)
  - Especially high dedupe with the corresponding client from other iterations – stress overlap affinity

- If new load outpaces capacity, high cost.  If the new appliance is added to keep up with added load, low cost.

**EMC²**

# When Capacity is Exceeded



Fractional Capacity Used

Clients

Cap w/o Dedupe ⊣△⊢    Cap w/Dedupe ⋅◇⋅

**EMC²**

# When Capacity is Exceeded



20 clients added per iteration

Add appliance

Fractional Capacity Used

Clients

Cap w/o Dedupe —△—     Cap w/Dedupe —◇—

# When Capacity is Exceeded



Offered load above 1.0 indicates overflow

20 clients added per iteration

Add appliance

Fractional Capacity Used

Clients

Cap w/o Dedupe — △ —    Cap w/Dedupe ⋯◇⋯

# When Capacity is Exceeded

# When Capacity is Exceeded

# When Capacity is Exceeded

# When Capacity is Exceeded

**EMC²**

# When Capacity is Exceeded



Legend:
- Cap w/o Dedupe △
- Cap w/Dedupe ◇
- Random ○
- Round Robin ✕
- Bin Packing +

(X-axis: Clients, 0 to 1000; Left Y-axis: Cost, $0$ to $10^6$; Right Y-axis: Fractional Capacity Used, 0 to ~1.5+)
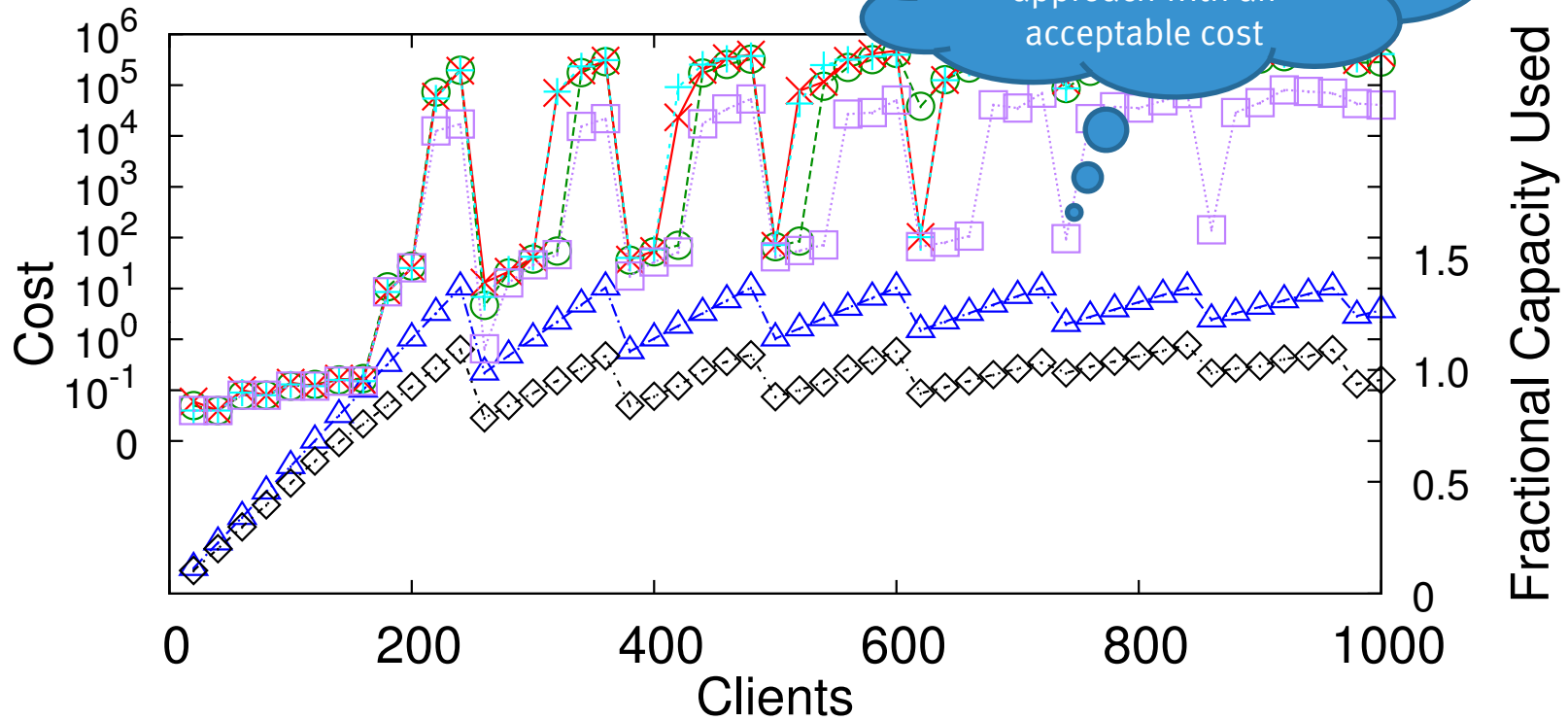
# When Capacity is Exceeded

# When Capacity is Exceeded

A few cases where annealing is the only approach with an acceptable cost



Legend:
- Cap w/o Dedupe △
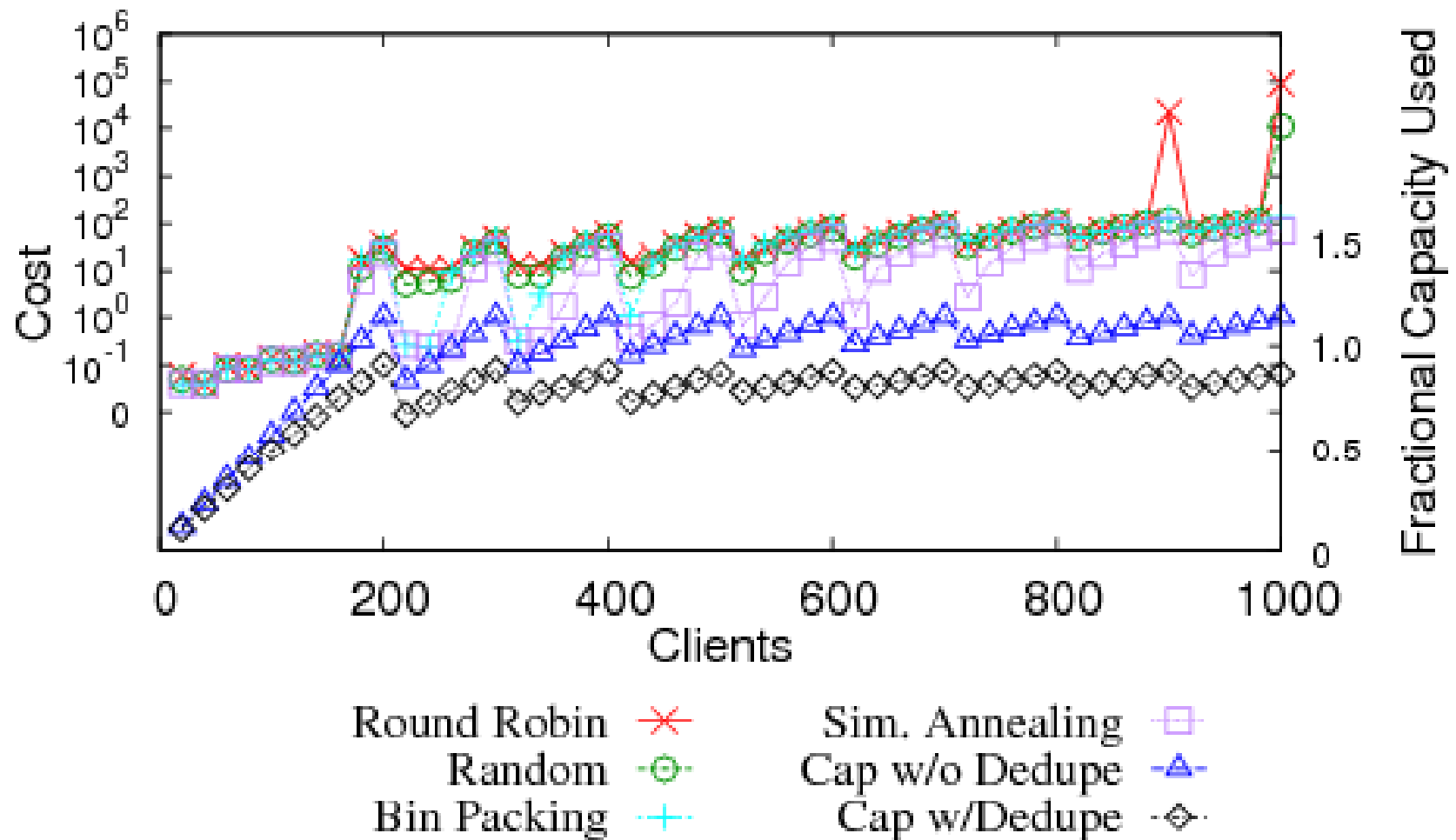- Cap w/Dedupe ◇
- Random ◯
- Round Robin ✕
- Bin Packing ✛
- Sim. Annealing ▢

EMC²

# When Capacity is Exceeded

# Roughly Fitting Within Capacity

# Roughly Fitting Within Capacity



Several cases where bin packing and annealing improve on the others when cost already low

Round Robin
Random
Bin Packing
Sim. Annealing
Cap w/o Dedupe
Cap w/Dedupe

# Roughly Fitting Within Capacity

Costs only occasionally very high



Legend:
- Round Robin — ✳ (red)
- Random — ⊙ (green)
- Bin Packing — + (cyan)
- Sim. Annealing — □ (purple)
- Cap w/o Dedupe — △ (blue)
- Cap w/Dedupe — ◇ (black)

# What Else?

- Refer to the paper for:
  - A more detailed discussion of overlap computation
  - Some other examples of using the assignment tool
  - Overhead analysis
    - Simulated annealing often works much better but is dramatically slower
  - Variants
    - Ignoring previous assignments
    - How to penalize for each client that doesn't fit
  - Impact of content-awareness

Backup slides for Q&A

EMC²

# Summary

- In a large IT environment, important to automate assignment of clients to backup appliances to optimize for <span style="color:red">capacity</span> and <span style="color:red">throughput</span>

- Taking content overlap into account can reduce capacity requirements and may improve throughput due to duplicate suppression

- Many options for how to balance load
  - All work well if not overloaded
  - Bin Packing somewhat better than the other simple techniques as limits approached
  - Simulated Annealing can handle some extra overload cases

EMC²

# THANK YOU