Competitive Clustering of Stochastic Communication Patterns on the Ring

Chen Avin Louis Cohen Stefan Schmid

Nice to meet you!

The Network Matters

Cloud-based applications generate significant network traffic
 E.g., scale-out databases, streaming, batch processing applications

E.g., Hadoop Terrasort job:



Virtual machine placement affects bandwidth costs



Virtual machine placement affects bandwidth costs



Virtual machine placement affects bandwidth costs



Virtual machine placement affects bandwidth costs





Option 1: Change the topology (?!)

Option 1: Change the topology (?!)

- □ Theory of demand-aware networks
- Prototypes emerging: e.g., ProjectToR (SIGCOMM 2016)
- Based on lasers and mirrors



Option 1: Change the topology (?!)

- □ Theory of demand-aware networks
- Prototypes emerging: e.g., ProjectToR (SIGCOMM 2016)
 - Based on lasers and mirrors

Option 2: Cluster the nodes

 Migrate frequently communicating nodes closer together

Option 1: Change the topology (?!)

- Theory of demand-aware networks
- Prototypes emerging: e.g., ProjectToR (SIGCOMM 2016)
- Based on lasers and mirrors



Option 1: Change the topology (?!)

- Theory of demand-aware networks
- Prototypes emerging: e.g., ProjectToR (SIGCOMM 2016)
- Based on lasers and mirrors

Option 2: Cluster the nodes

 Migrate frequently communicating nodes closer together

- □ Challenges of communication pattern clustering:
 - Communication patterns are not known ahead of time...
 - ... and may even change over time!

ion pa

Option 1: Change the topology (?!)

- ❑ Theory of demand-aware networks
- Prototypes emerging: e.g., ProjectToR (SIGCOMM 2016)
- Based on lasers and mirrors

Option 2: Cluster the nodes

Migrate frequently communicating nodes closer together

Challenges of commu

- Communication patter
- ... and may even

Thus: Need to repartition clusters in an online manner, depending on demand!









Example: 4 clusters of size 4



Now assume: changes in communication pattern!
 E.g., more communication (1,3),(3,4),(2,5) but less (5,6)

Example: 4 clusters of size 4



Now assume: changes in communication pattern!
 E.g., more communication (1,3),(3,4),(2,5) but less (5,6)



A simple and fundamental model (e.g., a rack):



ℓ servers ("clusters")

A simple and fundamental model (e.g., a rack):



A simple and fundamental model (e.g., a rack):





Problem inputs: k, ℓ , $\sigma = \{u_1, v_1\}, \{u_2, v_2\}, \{u_3, v_3\}, \dots$

Communication pattern over time

Problem inputs: k, ℓ , $\sigma = \{u_1, v_1\}, \{u_2, v_2\}, \{u_3, v_3\}, \dots$



Objective:

$$ALG(\sigma) = \sum_{t=1}^{|\sigma|} mig(\sigma_t; ALG) + com(\sigma_t; ALG)$$

Problem inputs: k, ℓ , $\sigma = \{u | Two flavors: (1) online (worst-case) pattern (2) learning: from a fixed (unkown) distribution$

Costs:



Objective:

$$ALG(\sigma) = \sum_{t=1}^{|\sigma|} mig(\sigma_t; ALG) + com(\sigma_t; ALG)$$

The Crux: Algorithmic Challenges

A) Serve remotely or migrate ("rent or buy")? When to migrate? If a communication pattern is short-lived, it may not be worthwhile to collocate the nodes: the migration cost cannot be amortized.

The Crux: Algorithmic Challenges

- A) Serve remotely or migrate ("rent or buy")? When to migrate? If a communication pattern is short-lived, it may not be worthwhile to collocate the nodes: the migration cost cannot be amortized.
- **B)** Where to migrate, and what? If nodes should be collocated, the question becomes where. Should the first node be migrated to the cluster of the second or vice versa? Or shall both be moved together to a new cluster? Moreover, an algorithm may be required to pro-actively migrate (resp. swap) additional nodes.

The Crux: Algorithmic Challenges

- A) Serve remotely or migrate ("rent or buy")? When to migrate? If a communication pattern is short-lived, it may not be worthwhile to collocate the nodes: the migration cost cannot be amortized.
- **B)** Where to migrate, and what? If nodes should be collocated, the question becomes where. Should the first node be migrated to the cluster of the second or vice versa? Or shall both be moved together to a new cluster? Moreover, an algorithm may be required to pro-actively migrate (resp. swap) additional nodes.
- **C)** Which nodes to evict? There may not exist sufficient space in the desired destination cluster. In this case, the algorithm needs to decide which nodes to evict, to free up space.

Online Variant: Competitive Ratio and Augmentation

Goal: minimize competitive ratio

$$ho(\mathrm{ON}) = \max_{\sigma} rac{\mathrm{ON}(\sigma)}{\mathrm{OFF}(\sigma)}$$

Online Variant: Competitive Ratio and Augmentation

Goal: minimize competitive ratio

$$\rho(\text{ON}) = \max_{\sigma} \frac{\text{ON}(\sigma)}{\text{OFF}(\sigma)}$$

□ Two flavors: without and with augmentation



Let's first look at special case: *k*=2







Special Cases: $\ell = 2$



Special Cases: $\ell = 2$





disk



- For 2 clusters: can emulate online caching!
 - k items, cache size k-1
- ❑ When item *i* is requested in original caching problem:
 - Introduce many requests between *d* and *i*: forces *i* to cache (if it is not yet)

<i>K</i> -1	d 00000	\longleftrightarrow	•
	00		

- For 2 clusters: can emulate online caching!
 - k items, cache size k-1
- ❑ When item *i* is requested in original caching problem:
 - Introduce many requests between *d* and *i*: forces *i* to cache (if it is not yet)
 - Which one to evict? Caching problem!

<i>k</i> -1		
	C	

- For 2 clusters: can emulate online caching!
 - k items, cache size k-1
- When item *i* is requested in original caching problem:
 - Introduce many requests between *d* and *i*: forces *i* to cache (if it is not yet)
 - Which one to evict? Caching problem!
 - Note: add many requests between *d* and nodes currently in cache: *d* stays in cache



disk

- For 2 clusters: can emulate online caching!
 - k items, cache size k-1
- When item is requested in original caching Lower bound k follows
 - Introduce many from caching! between a and i: 10 cache (if it is not yet)
 - Which one to evict? Caching problem!
 - Note: add many requests between *d* and nodes currently in cache: *d* stays in cache





Assume: requests only from a certain (ring) order



- Assume: requests only from a certain (ring) order
- Adversarial strategy: Whatever ON does, adversary will ask cut edge (exists even with augmentation): pays 1 each time!



- Assume: requests only from a certain (ring) order
- Adversarial strategy: Whatever ON does, adversary will ask cut edge (exists even with augmentation): pays 1 each time!
- Note: Adversarial request sequence only depends on ON! So online algo cannot learn anything about OFF.



- Assume: requests only from a certain (ring) order
- Adversarial strategy: Whatever ON does, adversary will ask cut edge (exists even with augmentation): pays 1 each time!
- Note: Adversarial request sequence only depends on ON! So online algo cannot learn anything about OFF.
- OFF can safely move to a partition which will be asked least frequently (once and forever)! Pigeon-hole principle: pays only every k-th time (i.e. k times less)



Online *Re*Partitioning: Overview of Results

- □ *k*=2 (online matching)
 - Greedy algorithm 7-competitive
 - Lower bound: 3-competitive
- O(k log k) -competitive algorithm CREP for 4-augmentation
 based on on growing components





Online *Re*Partitioning: Overview of Results



Learning Variant

- Adversary cannot choose request sequence but only the distribution
 - Adversary needs to sample i.i.d. from this distribution
 - Moreover: Adversary knows (deterministic or randomized) «learning» algorithm

Learning Variant

- Adversary cannot choose request sequence but only the distribution
 - Adversary needs to sample i.i.d. from this distribution
 - Moreover: Adversary knows (deterministic or randomized) «learning» algorithm
- Let's start simple: communication along ring only
 - I.e., adversary picks distribution over ring



Learning Variant

- Adversary cannot choose request sequence but only the distribution
 - Adversary needs to sample i.i.d. from this distribution
 - Moreover: Adversary knows (deterministic or randomized) «learning» algorithm
- Let's start simple: communication along ring only
 - I.e., adversary picks distribution over ring

Avoid high-weight edges on the cut!



J Naive idea 1: Take it easy and first learn distribution

- Do not move but just sample requests in the beginning: until exact distribution has been learned whp
- □ Then move to the best location for good

The Crux: Joint Optimi Learr

🕘 Naive idea 1: Take 🖌 🖌 asy an

Waiting can be very costly: maybe start configuration is very bad and others similarly good! Not competitive! Need to move early on, away from bad locations!

- Do not move but jussample requests in the beginning: until exact distribution has been learned whp
- Then move to the best location for good

J Naive idea 1: Take it easy and first learn distribution

- Do not move but just sample requests in the beginning: until exact distribution has been learned whp
- □ Then move to the best location for good
- Naive idea 2: Pro-actively always move to the lowest cost configuration seen so far

J Naive idea 1: Take it easy and first learn distribution

- Do not move but just sample requests in the beginning: until exact distribution has been learned whp
- ❑ Then move to the best location for good
- Naive idea 2: Pro-actively always move to the lowest cost configuration sector for



Bad: if requests are uniform at random, you should not move! Migration costs cannot be amortized. Crucial difference to classic distribution learning problems: guessing costs!



Mantra of our algorithm: Rotate!
 Rotate early, but not too early!
 And: rotate locally



Define conditions for configurations: if met, never go back to it (we can afford it w.h.p.: seen enough samples) Thm: Rotate ! Rotate early, but not too early! And: rotate locally



Mantra of our algorithm: Rotate!
 Rotate early, but not too early!
 And: rotate locally

If current configuration is **eliminated**, go to **nearby configuration** (in directed manner: no frequent back and forth)!



Mantra of our algorithm: Rotate!
 Rotate early, but not too early!
 And: rotate locally

If current configuration is **eliminated**, go to **nearby configuration** (in directed manner: no frequent back and forth)! **Growing radius** strategy: allow to move further only once amortized!





Conclusion

Dynamic repartitioning: a natural new problem!

- Competitive ratio super-linear in k: ok in practice (independent of number of servers!)
- Open questions:
 - Online variant: With less augmentation? Randomized?
 - Learning variant: General communication pattern, beyond ring?

