**Microsoft**

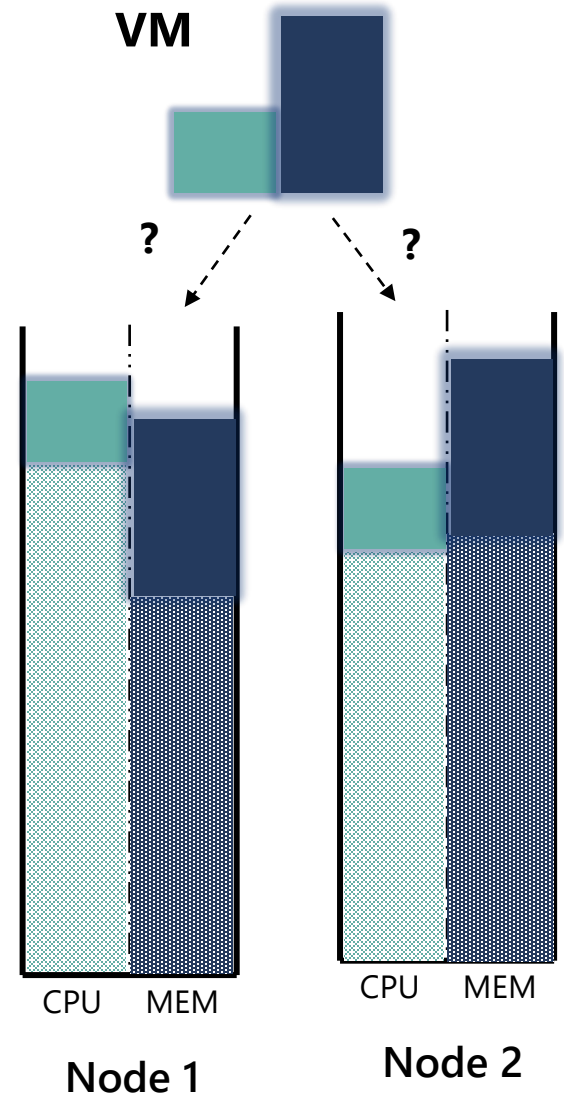# Virtual Machine Allocation with Lifetime Predictions

Hugo Barbalho, **Patricia Kovaleski**, Beibin Li, Luke Marshall, Marco Molinaro, Abhisek Pan, **Eli Cortez**, Matheus Leao, Harsh Patwari, Zuzu Tang, Tamires V. C. Santos, Larissa R. Gonçalves, David Dion, Thomas Moscibroda, **Ishai Menache**

# Motivation

- Allocation decisions have a direct impact on resource efficiency
- Inefficient placement might result in fragmentation and unnecessary over-provisioning
- Improvements of **1%** in packing efficiency can lead to cost savings of **hundreds of millions of dollars** (Hadary et al., 2020)

**Goal:** Increase Azure's packing efficiency with lifetime-aware algorithms

**Problem:** *Dynamic* multi-dimensional bin packing problem



**VM**

?  ?

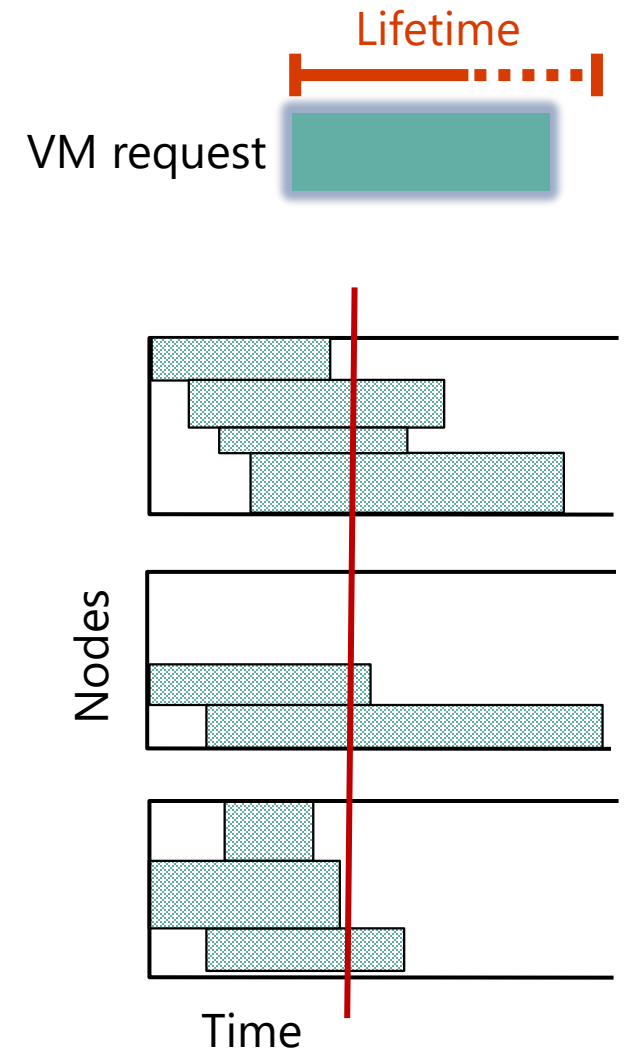CPU  MEM    CPU  MEM

**Node 1**          **Node 2**

# Motivation

- Allocation decisions have a direct impact on resource efficiency
- Inefficient placement might result in fragmentation and unnecessary over-provisioning
- Improvements of **1%** in packing efficiency can lead to cost savings of **hundreds of millions of dollars** (Hadary et al., 2020)

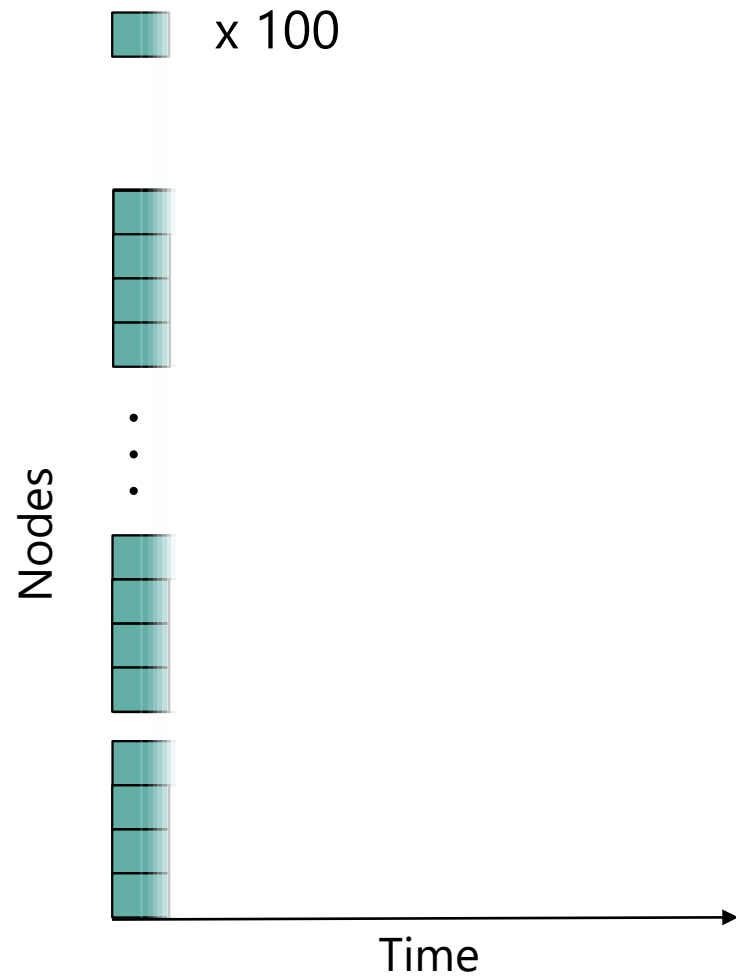**Goal:** Increase Azure's packing efficiency with lifetime-aware algorithms

**Problem:** *Dynamic* multi-dimensional bin packing problem
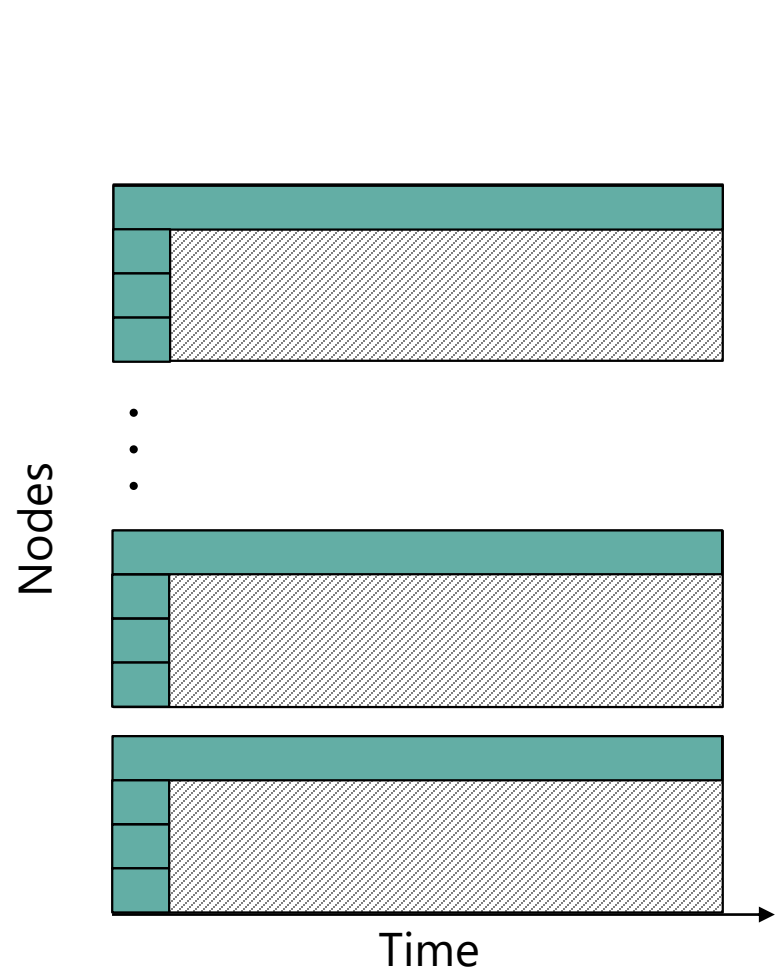
lifetime

Lifetime

VM request

Nodes

Time

# Example
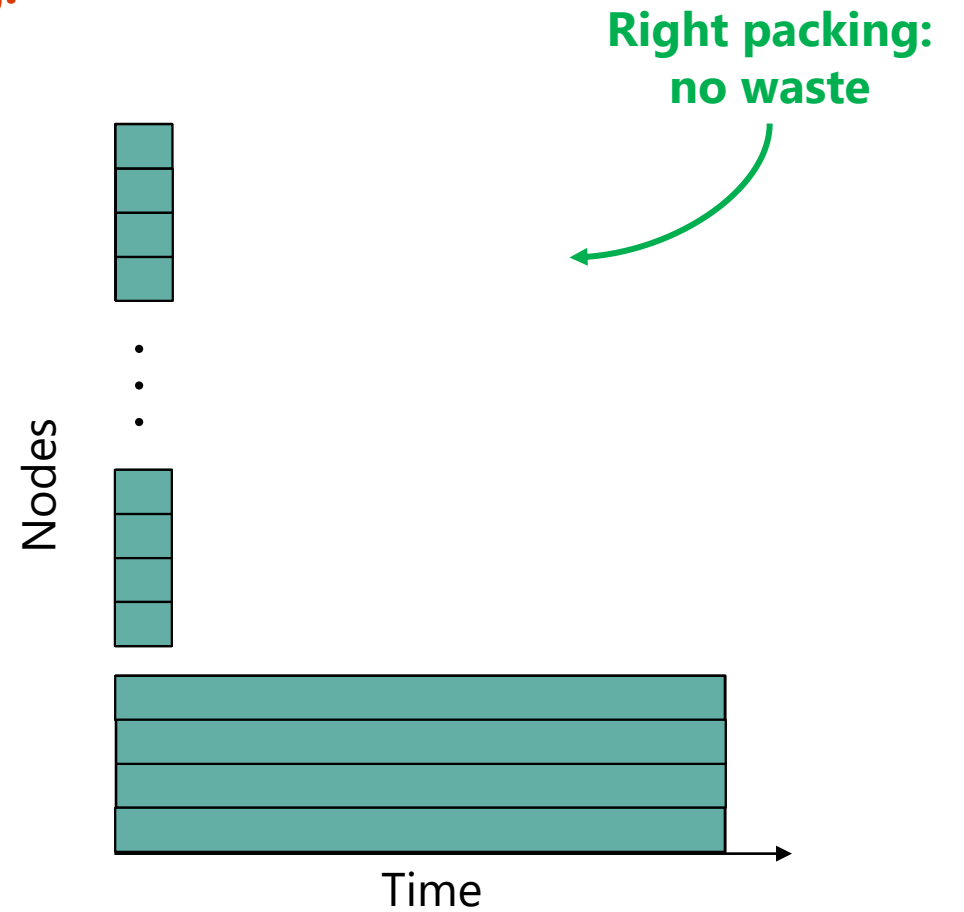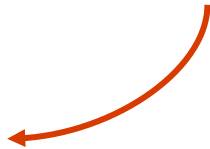
Why lifetime-aware allocations?
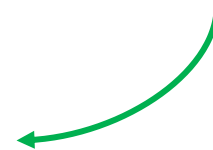
# Example

Why lifetime-aware allocations?

**Inefficient packing:
low density,
wasted resources**
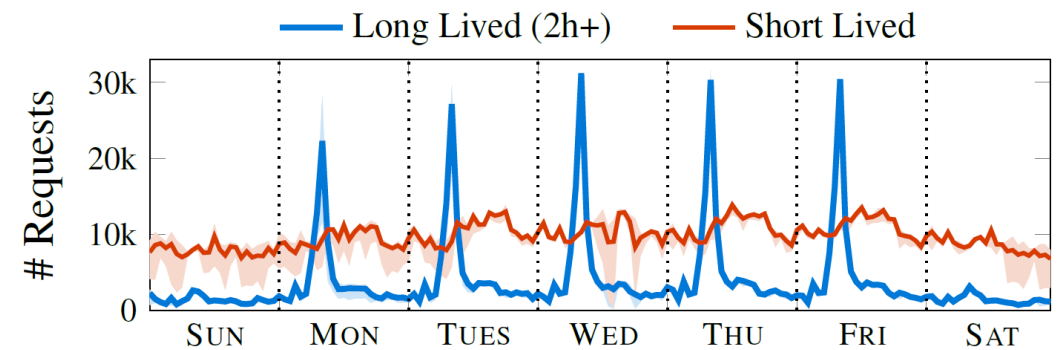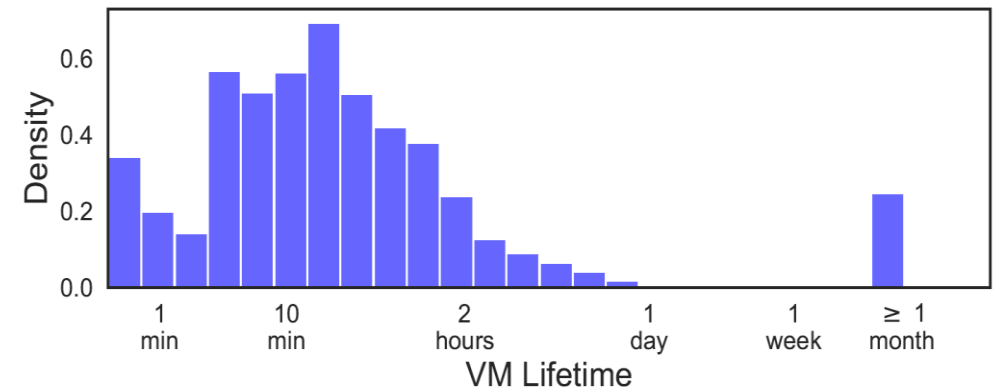
**Right packing:
no waste**

Nodes

Time

# VM lifetime characterization

- How are lifetimes in our system?

- High variance of lifetimes
  - Median: 16 minutes
  - Average: +1 day

- Lifetime temporal patterns
  - Feasibility of VM lifetime prediction

# Our contributions

1. Lifetime-aware **algorithm**
2. **ML model** for VM lifetime predictions
3. **System** to support it on real-time
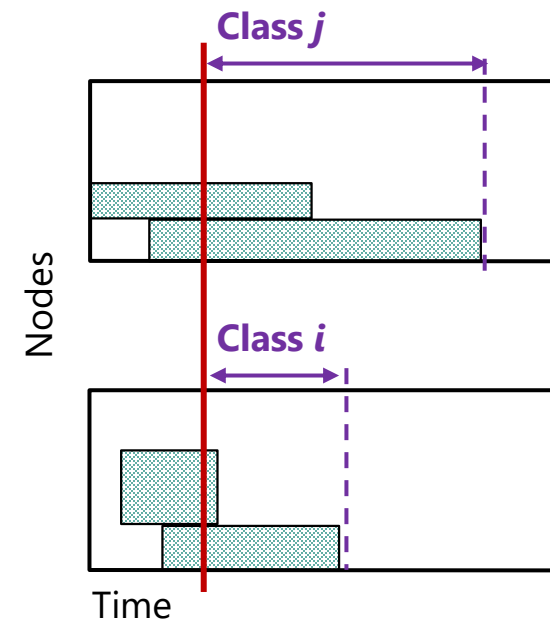
# Lifetime Alignment (LA) algorithm

**Idea:** "Prioritize putting jobs with similar lifetimes together"
· Lifetime ranges are partitioned into classes (where class 0 contains the smallest lifetimes)

For each incoming request:
- If the request is predicted **class 0**:
  - assign to **any** node using Best Fit

- If the request is predicted **class $j$**:
  - assign to a **class $j$** node (if exists), using Best Fit, else,
  - assign to **any** node using Best Fit

· Dynamically updates lifetime classification of nodes
  · Predicted remaining lifetime

· Theoretical indication that LA is **robust to prediction errors**
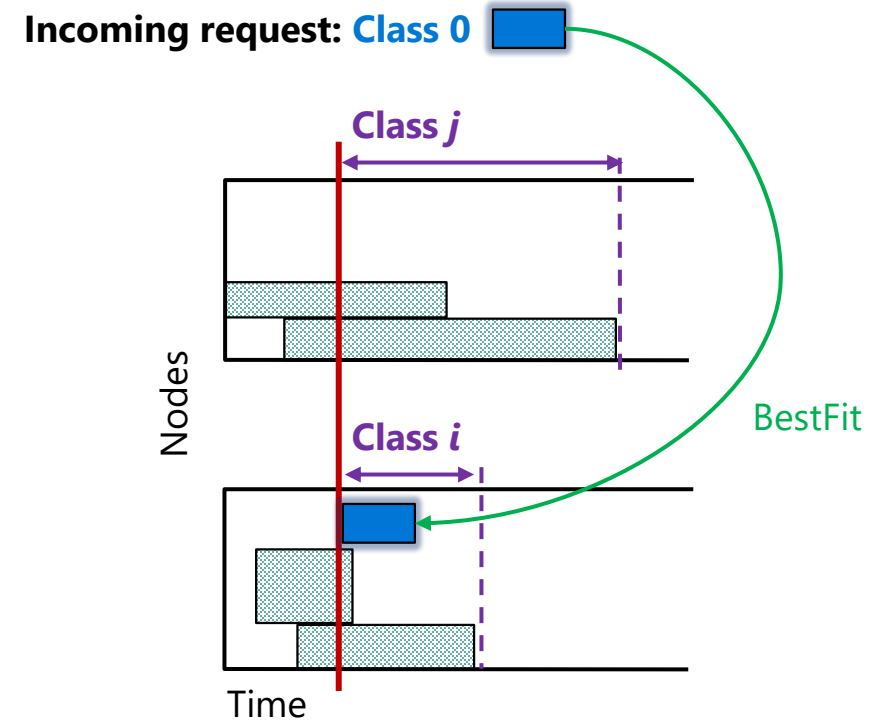
**Incoming request: Class 0**

Class $j$

Class $i$

Nodes

Time

# Lifetime Alignment (LA) algorithm

**Idea:** "Prioritize putting jobs with similar lifetimes together"

· Lifetime ranges are partitioned into classes (where class 0 contains the smallest lifetimes)

For each incoming request:
- If the request is predicted **class 0**:
  - assign to **any** node using Best Fit

- If the request is predicted **class _j_**:
  - assign to a **class _j_** node (if exists), using Best Fit, else,
  - assign to **any** node using Best Fit

· Dynamically updates lifetime classification of nodes
  · Predicted remaining lifetime

· Theoretical indication that LA is **robust to prediction errors**

# Lifetime Alignment (LA) algorithm

**Idea:** "Prioritize putting jobs with similar lifetimes together"

· Lifetime ranges are partitioned into classes (where class 0 contains the smallest lifetimes)

For each incoming request:
- If the request is predicted **class 0**:
  - assign to **any** node using Best Fit

- If the request is predicted **class $j$**:
  - assign to a **class $j$** node (if exists), using Best Fit, else,
  - assign to **any** node using Best Fit

· Dynamically updates lifetime classification of nodes
  · Predicted remaining lifetime

· Theoretical indication that LA is **robust to prediction errors**

**Incoming request: Class $j$**

Class $j$

Class $i$
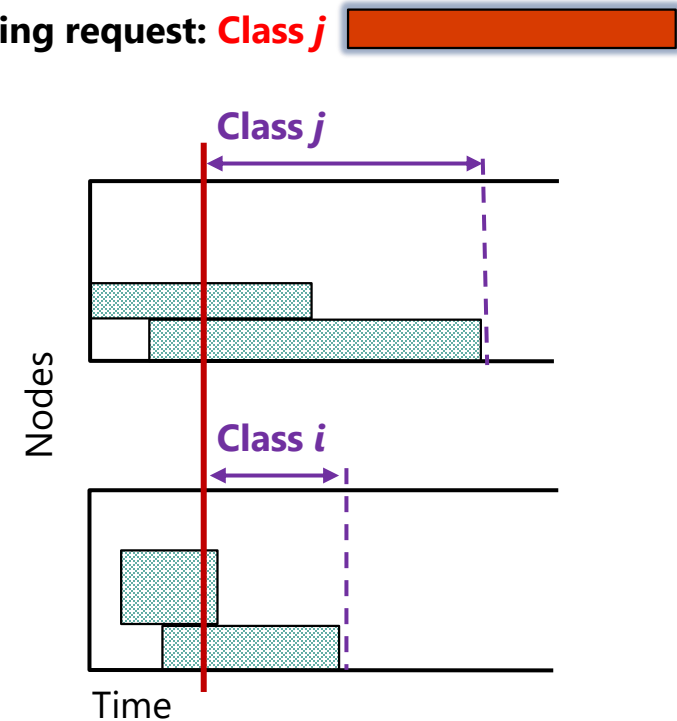
Nodes
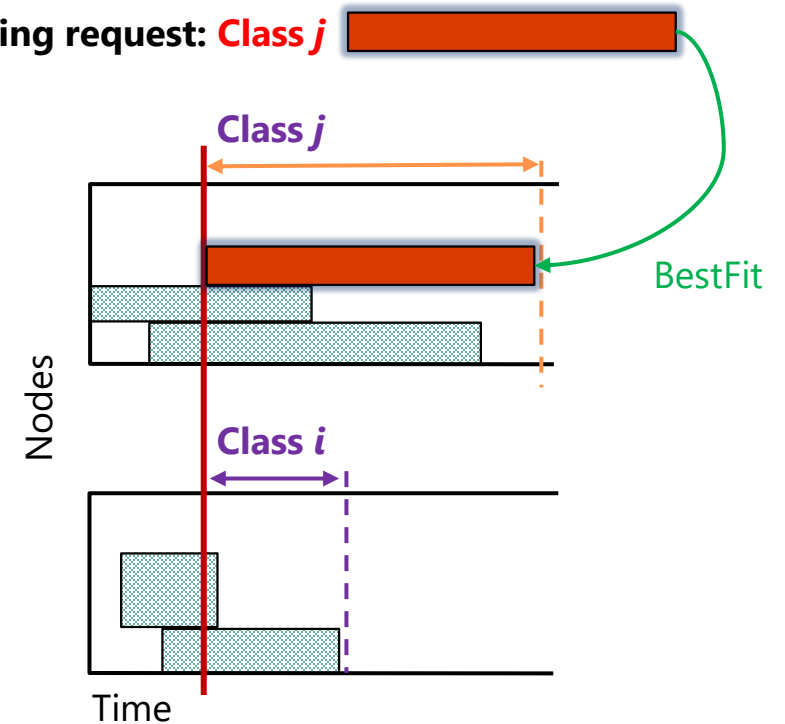
Time

# Lifetime Alignment (LA) algorithm

**Idea:** "Prioritize putting jobs with similar lifetimes together"

· Lifetime ranges are partitioned into classes (where class 0 contains the smallest lifetimes)

For each incoming request:
- If the request is predicted **class 0**:
  - assign to **any** node using Best Fit

- If the request is predicted **class *j*:**
  - assign to a **class *j*** node (if exists), using Best Fit, else,
  - assign to **any** node using Best Fit

· Dynamically updates lifetime classification of nodes
  · Predicted remaining lifetime

· Theoretical indication that LA is **robust to prediction errors**

**Incoming request: Class *j***

# Predicting lifetime

**Challenges:**

- Small feature set
- Fast inference time
- Missing data (loss or pruning)
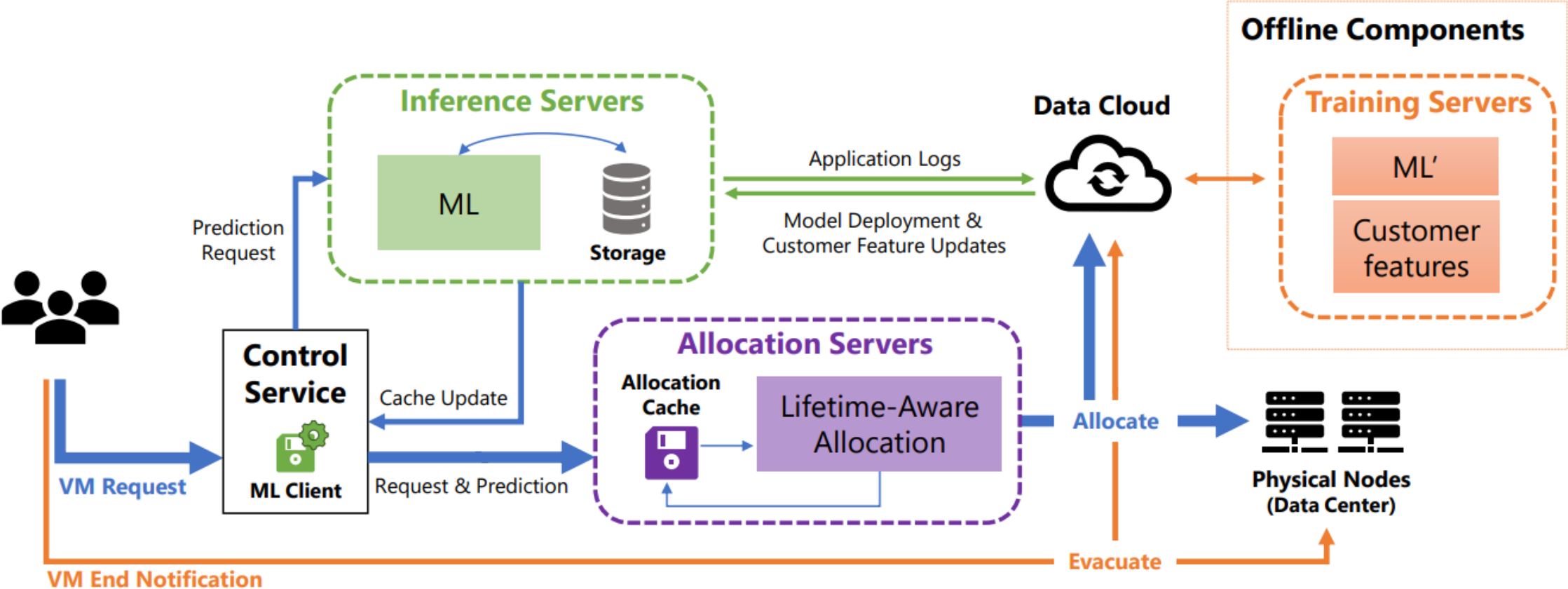- Skewed and long-tailed lifetime distribution

LightGBM model
Binary classification
Short/long threshold

Features:

- VM centric (VM type, OS, request time)
- Customer centric (temporal distribution)
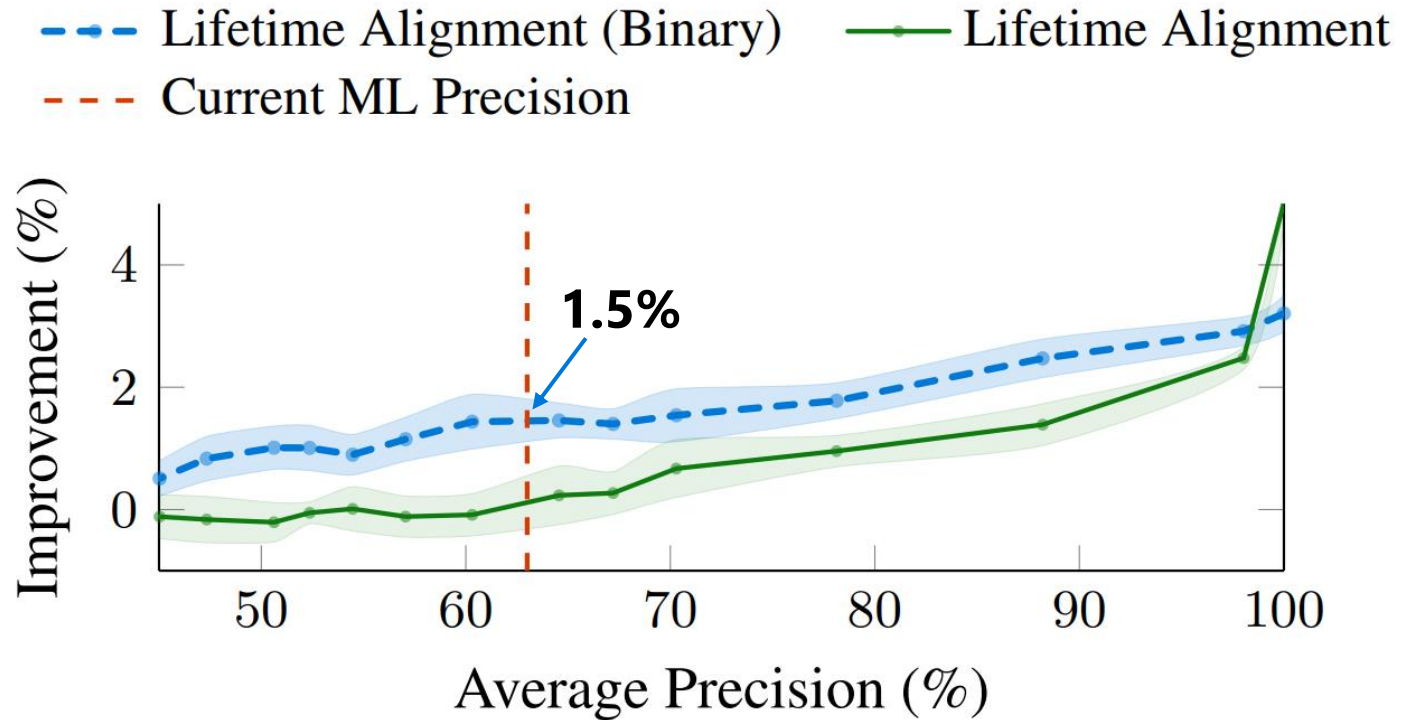
# System architecture

# Real-world production results

- Initial version (ML model + algorithm) in production

- 20 Million daily prediction requests
  - 200+ datacenters
  - 60+ regions

- 60% cache hit on inference results

- 99.2% predictions within time budget
  - Limit of 30ms

- ML model on production achieves expected performance

# Experiments



**Packing Density:** Measures the average number of allocated cores on non-empty machines

# Conclusion

We designed and implemented:

- **Lifetime-aware packing algorithm** robust to prediction errors
- **ML model** for VM lifetime predictions
- **System infrastructure** to support ML predictions in the critical path

➢ Packing improvements expected to save hundreds of millions of dollars per year

General methodology for resource management:

1. Produce data-driven intelligence (ML training, simulations) – offline, slower time-scale
2. Utilize the intelligence at real-time ("inference")
3. Applies to other scenarios, e.g., admission control (OSDI'23)

# References

Hadary, O., Marshall, L., Menache, I., Pan, A., Greeff, E. E., Dion, D., Dorminey, S., Joshi, S., Chen, Y., Russinovich, M., et al. Protean: VM Allocation Service at Scale. In 14th USENIX Symposium on Operating Systems Design and Implementation (OSDI 20), pp. 845–861, 2020.

Azar, Y. and Vainstein, D. Tight bounds for clairvoyant dynamic bin packing. ACM Trans. Parallel Comput., 6 (3), oct 2019. ISSN 2329-4949. doi: 10.1145/3364214.

Buchbinder, N., Fairstein, Y., Mellou, K., Menache, I., and Naor, J. Online virtual machine allocation with lifetime and load predictions. ACM SIGMETRICS Performance Evaluation Review, 49(1):9–10, 2021.

Microsoft

# Thank you