Cloud Spot Markets are Not Sustainable:

The Case for Transient Guarantees

Supreeth Subramanya, Amr Rizk, David Irwin UMassAmherst

Celling

Shared warehouse scale machines tend to have 10-50% utilization ⁹⁹

[**2013**] The Datacenter as a Computer: An Introduction to the Design of Warehouse-Scale Machines.



Commoditized compute



Users bid in a 2nd price auction

Idle Cloud Capacity





EC2 continually evaluates supplydemand to price spot servers



Allocate: bid price ≥ spot price Revoke: bid price < spot price



Commodity Spot Markets

Mature markets are inherently **VOLATILE**



Commodity and futures markets are great at pricing the resources and balancing supply and demand

but ...

Not possible to "beat the market" by **PREDICTING** future prices





Compute Time vs. Other Commodities





Compute time is **STATEFUL ***

∴ market volatility reduces amount the value of resources they allocate will decrease



Understanding Spot Market Characteristics





Spot Servers are Intrinsically Less Valuable!

Single-node batch job on a spot VM



•• On average, spot servers get less work done per unit of time compared to an equivalent on-demand server ••

Expected runtime





Spot Servers are Intrinsically Less Valuable!

Equilibrium Price of Spot (or price when spot stops being cheap)







• For this application, a spot server with **40%** discount on the on-demand price, provides no savings at all ⁹⁹





Distilling the Spot Market Characteristics



Available, Not Volatile, Predictable

Needs just one checkpointing

Needs as many checkpoints as there are revocations

We identify three key metrics: Availability, VOLATILITY, Predictability

Available, Volatile, Predictable

Available, Volatile, **Unpredictable**

Needs periodic checkpointing





Market Characteristics Impact the Performance





Equilibrium price of markets





On Spot Market Evolution







2009-2014



2015 onwards

Under mature market conditions

As they mature, cloud spot markets may not maximize the value of idle cloud capacity



11/15

• Uncertainty is more stressful than knowing for sure something bad will happen ⁹⁹

de Berker, Archy O., et al. "Computations of uncertainty mediate acute stress responses in humans." Nature communications 7 (2016)





Why Transient Guarantees?

Idle Cloud Capacity

EC2 Spot and GCE Preemptible

Transient Guarantees (MTTR based)



Not all spots are alike, and there are many ways to sell them

Transient Guarantees

E.g., Class-1 servers come with an MTTR of **55 hours**, and Class-4 servers **2 hours**

Able to value spot servers correctly Minimize fault-tolerance overhead

? Verifying transient guarantees

Providing probabilistic assurances on *availability, volatility and predictability* of spot servers

Increase revenue through differentiated offering Retain the freedom to reclaim any server

? Partitioning transient nodes into classes **?** Fixed pricing vs. market pricing

Thank you!

Supreeth Subramanya http://people.umass.edu/ssubramanya/

