



# Infinity Is Not a Strategy

## Right-Sizing the Cloud

SRECon-2026-Seattle

Praval Panwar

Principal Software Engineering Manager  
Observability @ Microsoft

# My Scaling-out Journey

.....



My house was my world

# Jodhpur



My city was my world....

\*courtesy govt of India



Country was my  
world



The planet  
earth was my  
known world

# Universe



Picture by NASA



Circinus Galaxy 13M light-years far

Hubble

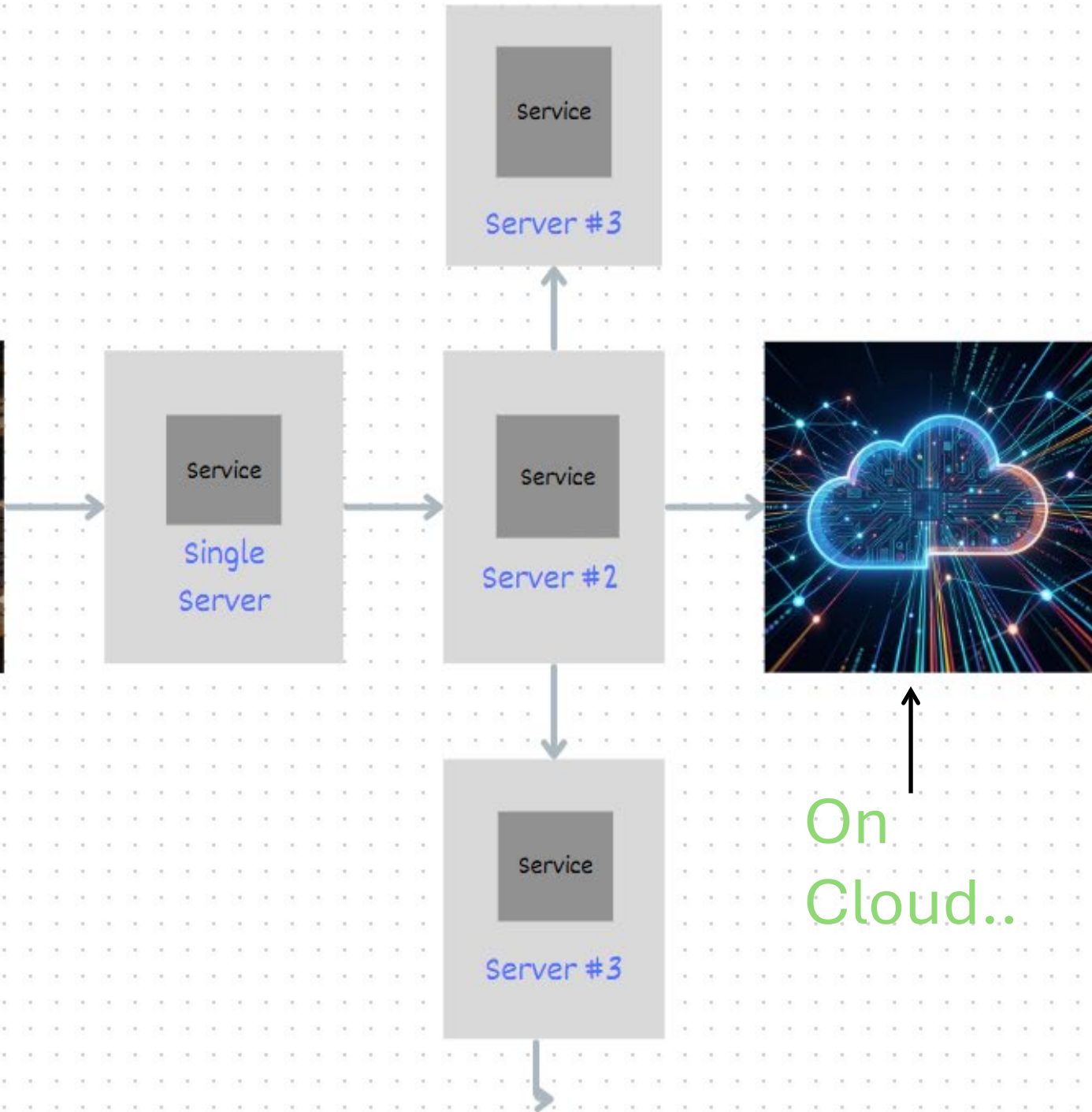
Picture by NASA

My own 'scale-out' capacity just kept growing as I grew ...

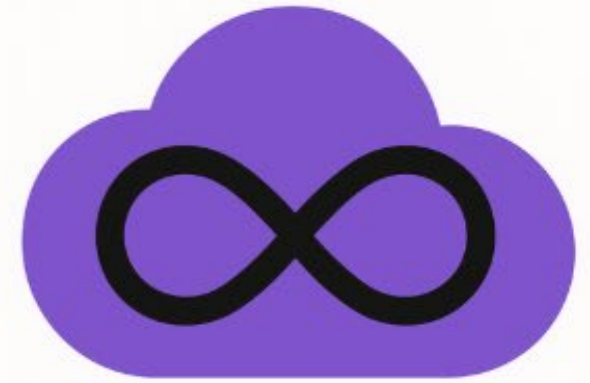
Drawing parallels from this example... ..



Hello  
World



On  
Cloud..



On a journey  
of going to  
infinite...

Scale-out can appear infinite at first glance, however.....

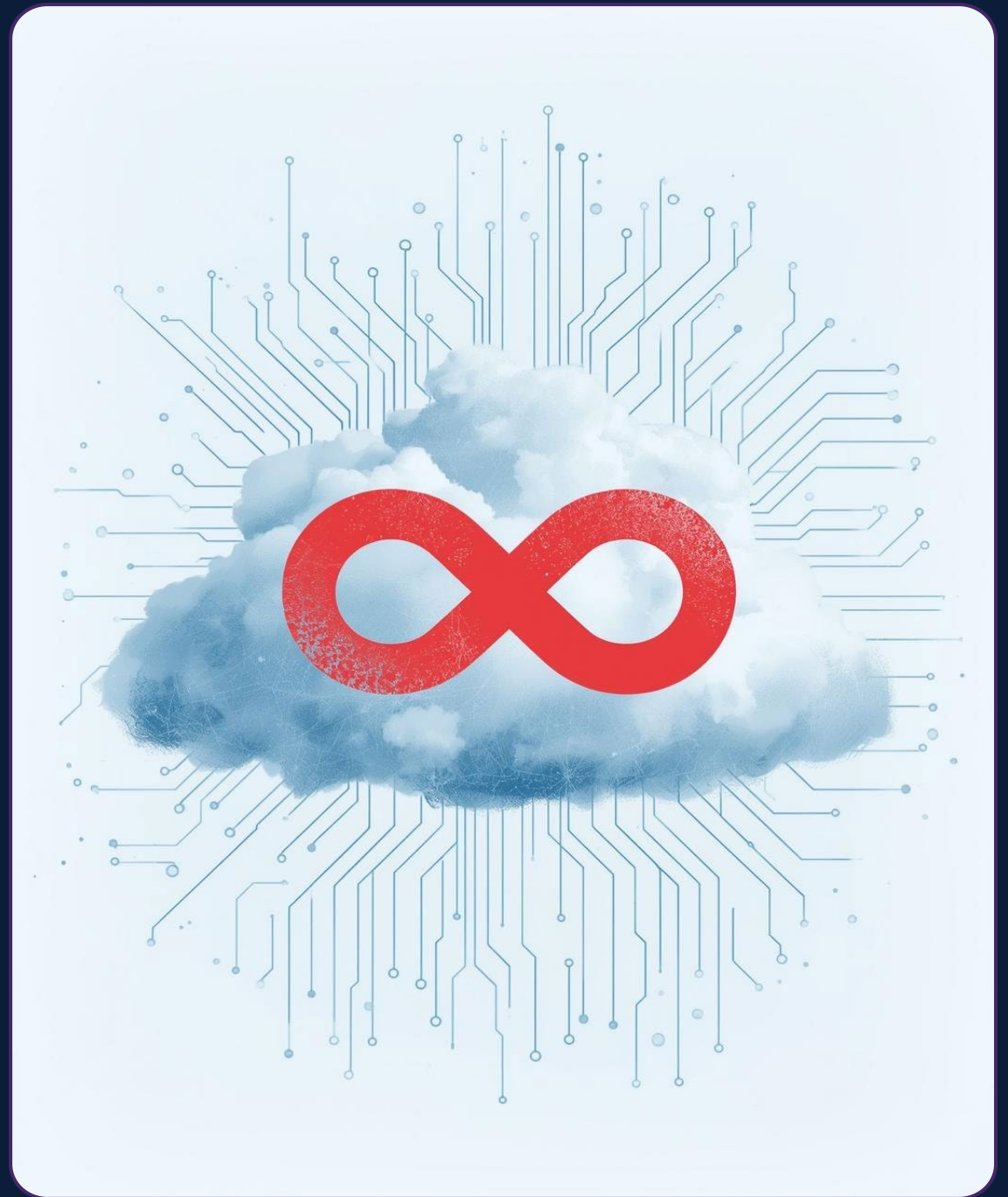


**Our finances  
are not infinite**

# The Fallacy of Infinite Scale

**The root problem is mindset, not tooling.**

---



- **‘The cloud is infinite’ is a comforting idea.**
- Every real system hits limits: quotas, noisy neighbors, budget, people...
- Can Foster a dangerous mindset. Often leads to insufficient planning and overspending.



**Datacenters.com** @datacenterscom · Jan 10



Cloud architecture used to be technical.  
Now it's a negotiation.

Capacity, pricing, power, and contracts decide what gets built—and where.

The cloud isn't infinite anymore.

#Cloud #CloudArchitecture #EnterpriseIT



**AITECH** @AITECHio · Jan 1

The Myth of Infinite Compute!

There is no such thing as infinite compute, only managed demand. Every AI system encounters trade-offs between speed, cost, and scale. What matters is whether those trade-offs are addressed deliberately or emerge unexpectedly.



**Micron21 Pty Ltd** @Micron21 · May 30, 2024



Organisations are noticing that leasing cloud resources has quite a high price tag, especially when compared to colocation. "Infinite" scalability is



**Open Intelligence Lounge**



**kasare ./** @lless\_tes · 2h

# the myth of infinite cloud scale

everyone assumes cloud scales forever  
more racks, more power, more cooling  
and somewhere, more intelligence  
but physics does not negotiate



**WindowsForum** @windowsforum · Aug 15, 2025



Surprise! "infinite" cloud capacity hit a limit in East US, leaving admins in a virtual jam. Turns out, even clouds have their rainy days!

**Forbes**

Many companies continue to overprovision cloud resources, paying for capacity that is not being used. According to the Azul survey, nearly 70% of companies say they are paying for cloud capacity that they are not using, and more than 40% say they use less than 60% of the public cloud compute they are paying for.

- **The message is consistent: infinite scale is a myth.**
- **Capacity is finite, physics is real, and most of us are massively overprovisioned.**
- **This illusion of infinite scale also comes with a hefty bill believe it or not...**

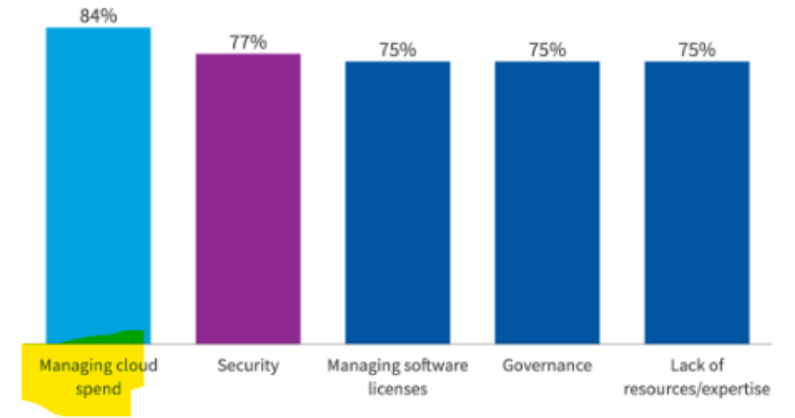
# What Companies Can Do About Cloud Spend Wastage

By [Eldar Tuvey](#), Former Forbes Councils Member.  
for [Forbes Business Council](#), COUNCIL POST | Membership (fee-based)

Oct 8, 2024 9:00 AM Eastern Daylight Time

## New Survey Finds Cloud Waste is On the Rise - Driven by Preventable Mistakes, Inefficiencies, and New AI Initiatives

### Top cloud challenges for all respondents



N=759  
Source: Flexera 2025 State of the Cloud Report

**flexera**

## Companies flush money down the drain with overfed Kubernetes cloud clusters

Just 13% of provisioned CPUs, 20% of memory utilized, study finds

 [Dan Robinson](#)

Fri 1 Mar 2024 | 09:30 UTC



A nontrivial  
problem

(sometimes)...



# **World Cup Tickets**

# FIFA World Cup 2026 ticket frenzy unfolds

- What happened wasn't necessarily a bug. It was a capacity problem driven by uncertainty.
- When tickets open, demand doesn't arrive gradually. It arrives as a sudden burst.
- Maybe it's two million users. Maybe ten million. You don't know in advance... You can try to guess

- You provision for average demand; but the system becomes overloaded during these bursts
- This is what capacity planning under uncertainty looks like in the real world
- Demand is unpredictable, and your job isn't to provision infinite capacity
- ✓ It's to provision **intelligently**.

Capacity  
Planning  
Under  
Uncertainty

---

Toilet paper and  
Covid

# The psychology behind why toilet paper, of all things, is the latest coronavirus panic buy

By [Scottie Andrew](#), CNN

🕒 5 min read · Updated 5:14 PM EDT Mon March 9 2020

- Production capacity hadn't suddenly collapsed.
- The infrastructure was still there.
- What changed was demand behavior.

- The system wasn't broken. It was operating exactly as designed
- *Just outside its expected demand range*
- Limits of known and limits of unknown

- SRE -> 100% reliability does not exist.
- Systems operate in a permanently degraded mode..
- 100% reliability = *Infinite redundancy*. -> *Infinite failover* -> *Infinite capacity*.
- It's impossible to achieve.

- Capacity planning operates under the exact same constraint.
- You will never have 100% certainty about the load your system will receive...
- The goal is to provision *intelligently*—to operate efficiently within the known range, and degrade gracefully outside of it



# Resource Misallocation

When faced with uncertainty, there are two instinctive ways to design systems.

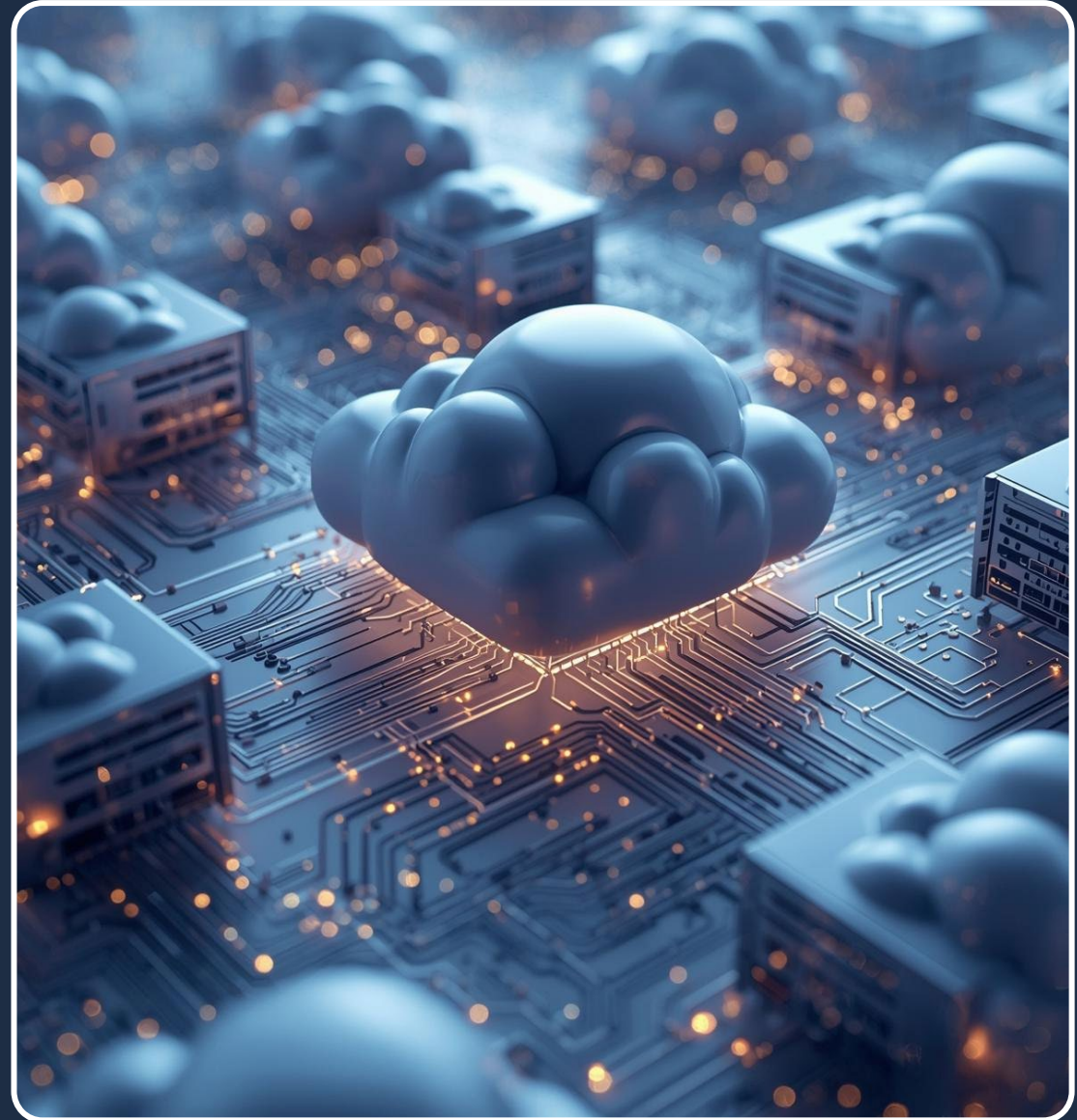
- If you build a fortress, you misallocate resources through over-provisioning.
- If you walk the tightrope too closely, you misallocate risk—and the system becomes fragile.
- Understand service fragility and operate within acceptable risk boundaries.

# How to balance Risk and Capacity



# Right-Sizing Explained

**Strategic adjustment** of cloud resources to align with actual workload needs in cost-effective manner



- The ideal scenario is reaching a right-sized state, when resources are optimized for workload demands
- While balancing cost efficiency with high performance to ensure seamless operations
- It's not an absolute number; it changes depending on system demands and you may need to have a dynamic approach to it

# Cross Industry Playbook for SRE

---





**Every seat, every  
flight, \$99 always.**

How do you  
compete with that?



# American Airlines: Turning Data Into Strategy: SABRE

How many \$99 seats should we sell to fill the plane, and how many should we hold back for business travelers who'll pay \$300 the day before.

data turned guesswork into a calculation.

# How American Airlines Leveraged data

The history of SABRE, yield management, and modern airline economy we know it.

*American could price tickets such that every seat on every flight was bringing in the maximum dollar amount that each customer was willing to pay, giving them the flexibility to take down People Express.*

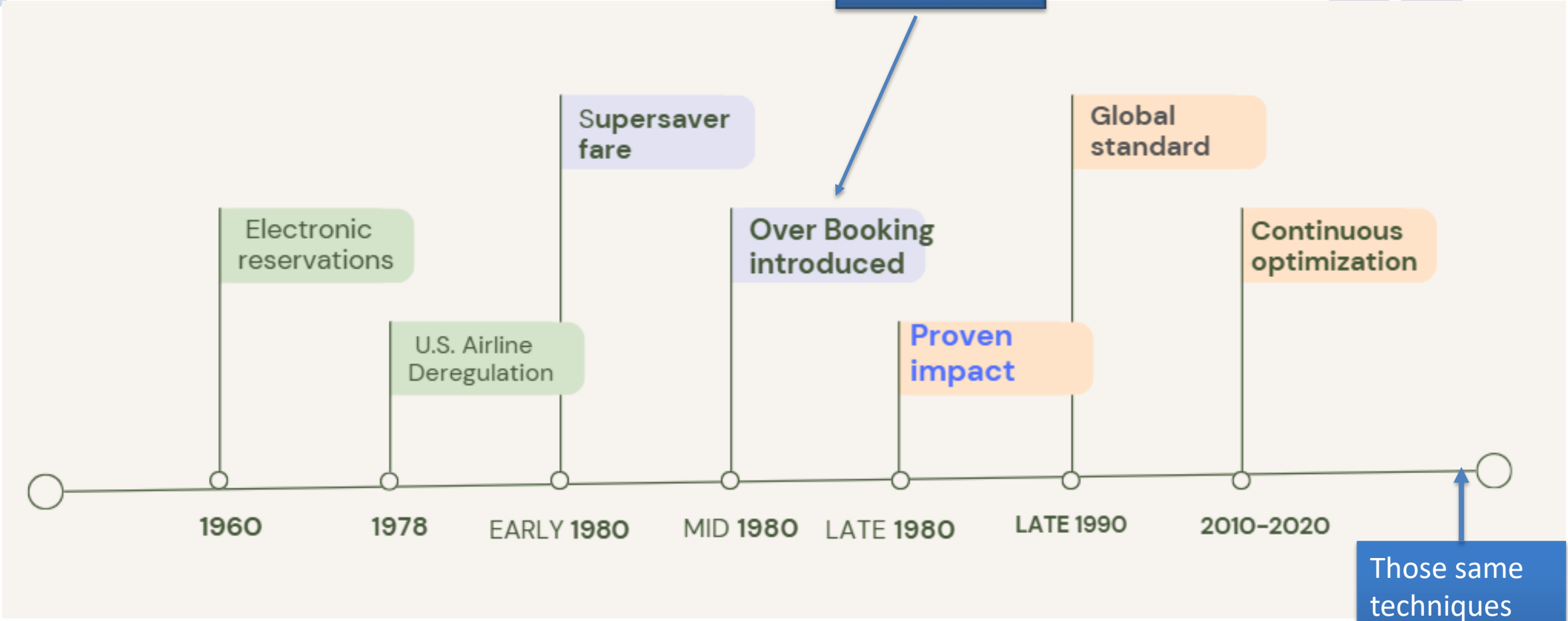
American Airlines deployed a system that Used historical booking data to forecast demand and price seats dynamically

They modeled no-show rates per route, per time, per day of the week. If expect 10 no-shows, the sell 10 extra tickets.

advantage of the fact that many customers either cancel their airline reservations, or do not show up for their flights at all (cf. Suzuki, 2006). American Airlines reported in 1990 that the benefit of overbooking to the airline exceeded \$225 million (cf. Smith et al.,



ARIMA,  
Booking curve  
and  
overbooking



Those same  
techniques  
are still in use  
today

# Airlines: ARIMA

$$\text{Demand}_t = c + \sum \phi_i \cdot \text{Demand}_{t-i} + \sum \theta_j \cdot \epsilon_{t-j}$$

- Problem : How many people will show-up for this flight?  
Answer : use ARIMA -> forecasts passengers per flight.
- These models learn patterns over time: trend, seasonality, weekday/weekend effects.
- They give you a forecast and often an uncertainty range, not just a single number

# Airlines: Booking Curves

- Booking curves: plot the fraction of seats sold versus days before departure
- For given route -> curve of 20% of seats sold 60 days out, 60% by 14 days out, 90% the day before
- Compare the live booking curve to historical curves and adjust the forecast
- **Action: Use time-series models to forecast traffic per service + use “booking curves” to forecast launch traffic**

# Airlines: Overbooking $\rightarrow$ Explicit Overcommit

Total Seats booked =  $C + k$ ;  $P(X < k)$

$\Rightarrow$   $C$  = physical seats,  $k$  = overbooking constant

$\Rightarrow$   $X$  = random variable for no show

- What's the distribution of passengers who don't show up
- Pick  $K$  such that more tickets sold  $\Rightarrow$  more revenue
- Denied boardings become too common or costly, then back off
- Probabilistic Overcommitment: A controlled, data-driven risk trade-off, and not accidental

# Airlines: Overbooking → Explicit Overcommit

- $C$  is physical CPU/memory on a host;  $k$  is the extra ‘virtual’ capacity we allocate via overcommit
- *‘What is the probability that all workloads spike together?’*
- Set  $k$  to keep overload probability below a chosen threshold

**Action: Pick overcommit factors from utilization distributions and define risk/SLO**

You overcommitted, what  
if you don't have system  
for tracking uncertainty

No system to track  
overprovision or real-time  
control?



BUSINESS

## In JetBlue's wake, a push for fliers' rights

By Martin Zimmerman, Alana Semuels and Molly Selvin

Feb. 20, 2007 12 AM PT



TIMES STAFF WRITERS

Thousands of passengers were stranded in New York. Flights were scrubbed at airports around the eastern U.S. But the biggest fallout from JetBlue Airways' monumental meltdown over the last week may land on Capitol Hill.

The travails of the nation's eighth-largest carrier might not be enough on their own to prompt federal action on behalf of air travelers. But JetBlue's troubles come as the airline industry struggles through a winter of discontent punctuated by airport shutdowns, massive flight cancellations and tales of passengers trapped for hours aboard parked aircraft.

BUSINESS | JetBlue's C.E.O. Is 'Mortified' After Fliers Are Stranded

The basic problem, he said, was JetBlue's communication system: the ice storm had left a large portion of the airline's 11,000 pilots and flight attendants far from where they needed to be to operate the planes, and JetBlue lacked the trained staff to find them and tell them where to go. Prior to last week, JetBlue had never had so many people out of position.

The reservation system was also overwhelmed, with customers

- 156 flights scheduled that day. Only 17 flew
- 11,000 pilots and flight attendants out of position. They had no way to communicate
- Created real-time crew tracking database
- This is why crew rotation optimization exists in the industry.

# Airlines: N/W, Rotations, → Bin-Packing & Failure Domains

- Mixed-integer programming problem
- Flights form a graph; Planes/crews are flows with constraints
- If disruption happens; you need to re-optimize in real time; with minimum-cost recovery plan

# Airlines: N/W, Rotations, → Bin-Packing

Services -> flights; Compute resources -> crews

Regions and availability zones -> airports

- When a region goes down -> do we have a plan to re-route traffic and workloads in real time?

**Action: Schedule workloads with anti-affinity; decide what to cancel first—like low-priority batch work**

# Airlines → Cloud Capacity: Key Lessons

- ✓ Forecast with structure like ARIMA  
+ booking curves
- ✓ Overcommit explicitly, not accidentally
- ✓ Rerouting load with clear priorities

All of these push us away from 'infinite scale' thinking....

# Power Grids: Real-Time Balance





# Power of a Tree

---



BLOG > A LOOK BACK AT THE NORTHEAST BLACK...

## A Look Back at the Northeast Blackout of 2003 and Lessons Learned

August 10, 2023



**In August of 2003, the largest blackout in North America occurred, affecting 50 million people at an estimated cost of \$4 - \$10 billion**

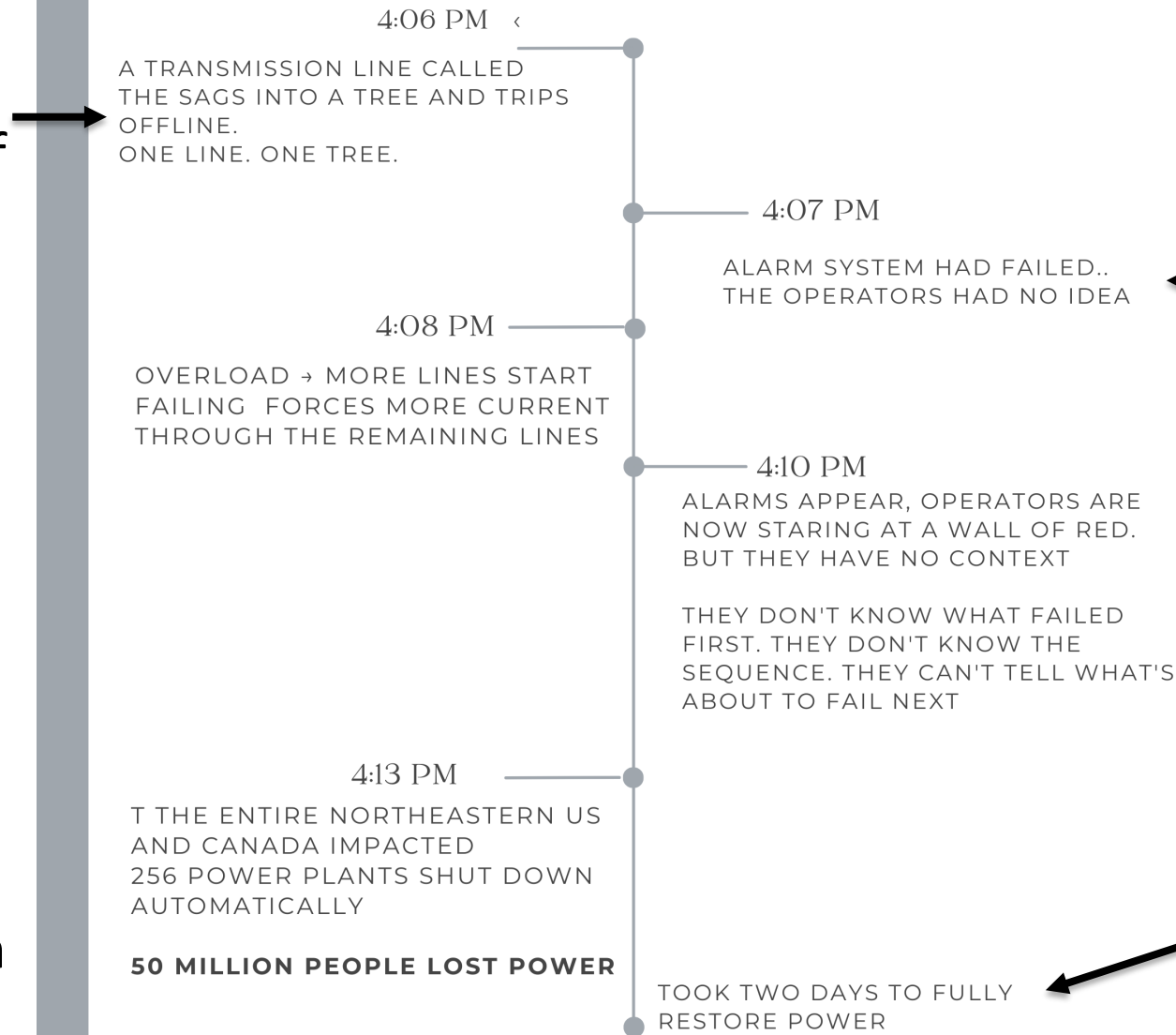
### Proximate Causes:

- Load imbalance caused by generator shutdown triggered cascading transmission line failure

On August 14, 2003, a transmission line fault in Ohio caused by contact with a tree cascaded into what would become one of the largest outages in North American history plunging more than 50 million people in eight states and Ontario into darkness.

# AUG 14TH, 2003

Power grids are designed for this, the whole point of having a **network..**



The software bug was silent

New York City, Cleveland, Detroit, Toronto, Newark  
Airports/trains shut down

Conclusion  
**The grid violated N-1**

**N-1 is the principle that says: your system must survive the loss of any single major component.**

One generator. One transmission line.  
One substation.

If losing one thing causes everything else to collapse, you don't have a resilient system—you have a house of cards.

Failure is often the trigger for innovation  
The *industry rebuilt*.

- ✓ They mandated real-time monitoring systems that don't fail silently.
- ✓ Requiring to run N-1 contingency analyses
- ✓ They formalized Optimal Power Flow algorithms to keep the grid operating inside safe limits and automatic load-shedding schemes.

# Power Grids: N-1 & Reserve Margin

---

- **N-1**: we must meet SLOs if any one AZ/region/datastore fails
- +
- **Reserve margin**: plan capacity above forecast peak
- +
- **Spinning reserve**: run some clusters below their max, ready to ramp up immediately

**Action : Define N-1 (AZ/region) dependency and pick headroom from failure simulations + pre-defined load shedding for low-priority work**

NATIONAL

It's been five years since **catastrophic**  
**Texas blackouts**. How much has

FEBRUARY 13, 2026 · 4:34 PM ET

HEARD ON ALL THINGS CONSIDERED

By Natalie Weber

DEITABASE

Home

About

Case Summaries

## 2021 Texas Power Grid Failure – a preventable disaster

### Description and Details

The failure of Texas' power grid in February, 2021 was one of the most severe energy crises in U.S. history, leaving millions of people in freezing temperatures. The roots of this disaster lie in Texas' unique power grid and approach to energy regulation.

- Feb 15, 2021: Winter Storm Uri hits Texas
- Planned for 67 GW winter peak demand per previous year
- Actual demand: 76 GW (9 GW over forecast)
- 40%+ of generation (30k MW) went offline simultaneously
- Grid came within 4 minutes 37 seconds of total collapse

# Power Grids: Short-Term Load Forecasting

---

- ✓ That nine-gigawatt gap—is exactly what **short-term load forecasting** is supposed to close
- ✓ If your capacity planning is based on '*last season peak times 1.3*' and you never update the model, you will get surprised...
- ✓ Grids now use **continuous load forecasting models** with real-time weather data, neural networks.

# Power Grids: Short-Term Load Forecasting

---

- **Inputs:** history, weather forecasts, calendar (day/holiday)

$$\hat{L}_{t+1:t+H} = f(L_{t-k:t}, \text{weather}_{t:t+H}, \text{calendar})$$

forecast of the past  $L$  values, where  $f$  is the learned neural network.

- **Models:** a neural network—often an LSTM or hybrid
- **Output:** predicted load over the next hours or days, with an error distribution

# Power Grids: Short-Term Load Forecasting

---

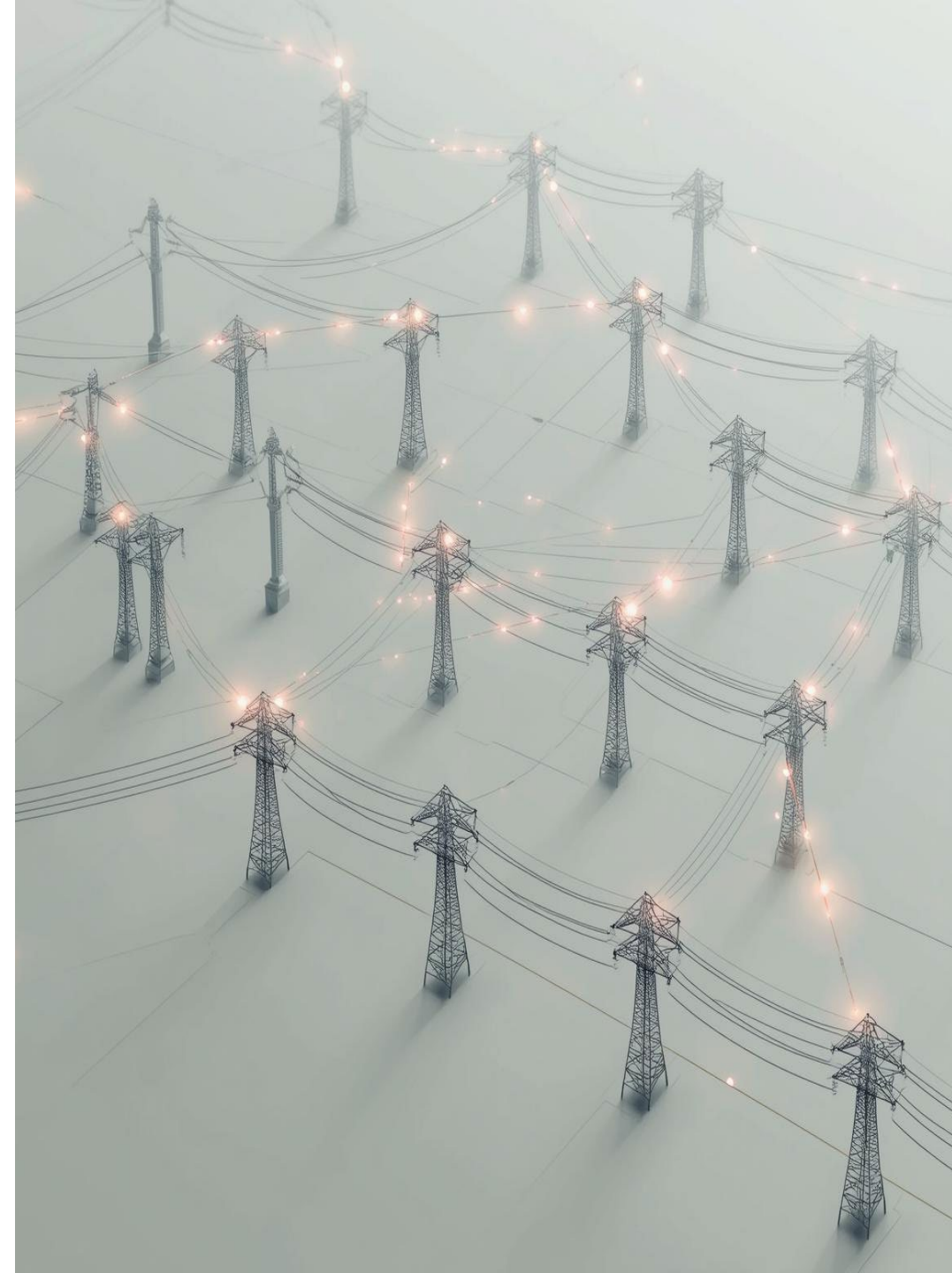
- Replace **L** with your QPS or CPU utilization per service/region
- Replace 'weather' with external drivers: marketing campaigns, product launches, seasonality
- Even a simple model—like Prophet or an LSTM—can be your  $f()$  and give you better forecasts

**Action: Treat load forecasting as a first-class discipline with continuous model updates, feedback and real-time inputs**

# Power Grids → Cloud: Key Lessons

- ✓ Explicit N-1/N-2 targets
- ✓ Conscious reserve margins
- ✓ Predefined load shedding

*Treat capacity as risk  
management, not infinite  
elasticity*



# Logistics: Moving Atoms at scale



---

Remember  
When Toilet  
Paper Became  
a Luxury Item?

---



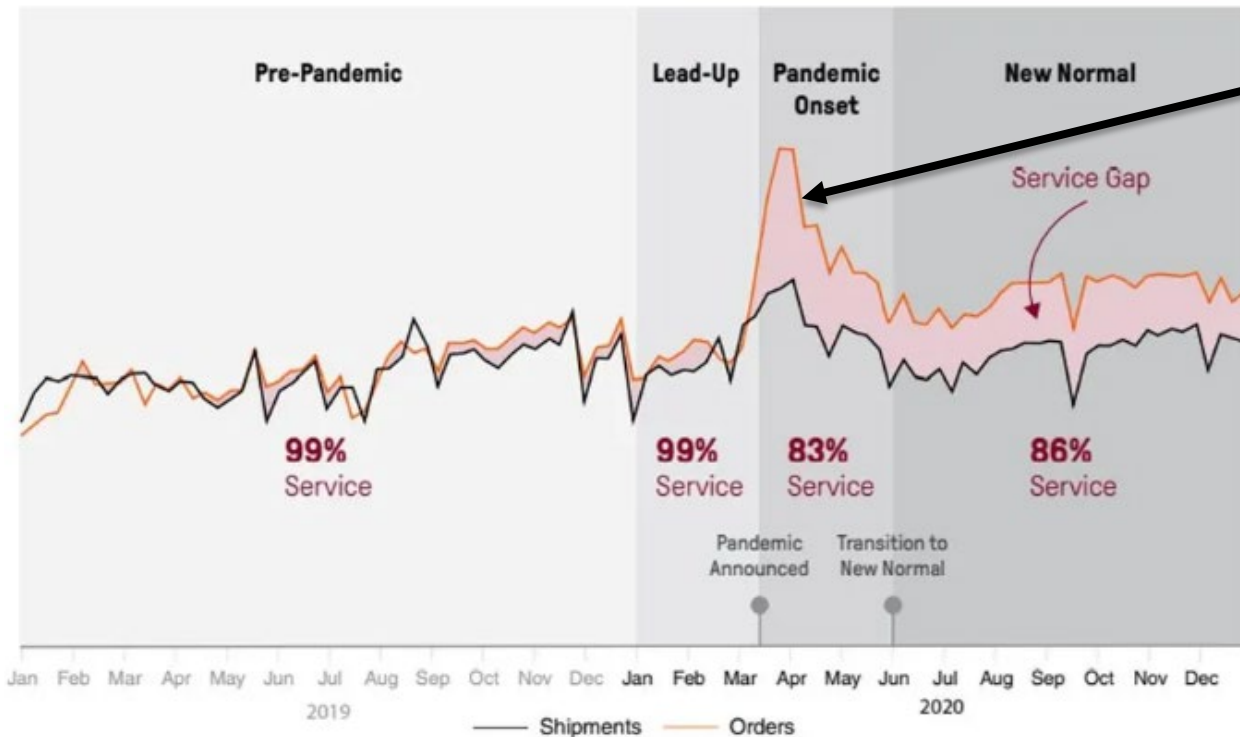
- Traditional forecasting : What sold last March, what sold the March before that—and you use those patterns to predict what you'll need this March
- Add some seasonal adjustments, some trend smoothing, and you've got your forecast. It works. *Most of the time.*
- **But March 2020 was not like any previous March.**

# Toilet Paper Shortages, Empty Shelves, And Panic Buying: Just How Bad Was Grocery Service In 2020?

By [Steve Banker](#), Contributor. © I cover logistics and supply chain management.

[Follow Author](#)

Published Oct 01, 2021, 09:23am EDT, Updated Oct 01, 2021, 01:22pm EDT



Demand has **exploded**.

- Households are using 40% more toilet paper than normal.
- Demand planning error increased to 59%

Forecasted v/s Sold



Rapid demand response forecasting helps retailers adapt during COVID-19

by **CHARLIE CHASE** on JUNE 11, 2020

0 COMMENTS

Reduced forecast error from  
20-60% down to 4-9%.

Some companies did better,  
The study found that orgs  
using demand sensing (*real-time signals – sales, search*)  
experienced fewer errors....

Do you know what was the  
difference?

✓ **Exponential Smoothing**

Gives more weight to recent  
data than older data

# Logistics: Exponential Smoothing

$$\hat{D}_t = \alpha \cdot D_t + (1 - \alpha) \cdot \hat{D}_{t-1}$$

$D_t$  = actual demand at time  $t$ .

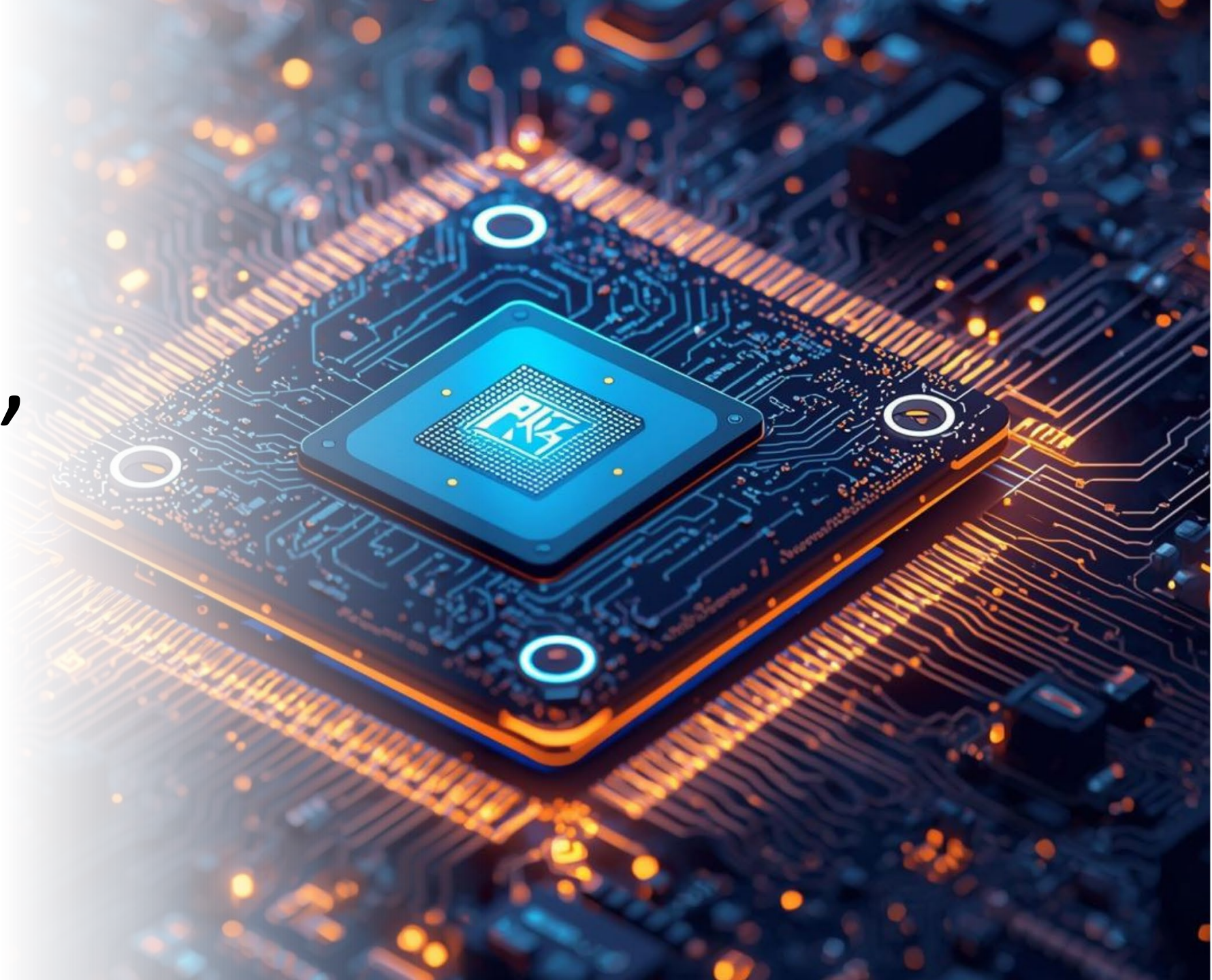
$\hat{D}_t$  = updated forecast.

$\alpha \in (0,1)$  controls how fast the forecast reacts.

- **Recent data gets more weight; older data decays exponentially**
- **Better stocking decisions when combined with variability metrics**

**Action: Use exponential smoothing for internal services' capacity baselines or autoscaling target**

No Chips,  
No Cars



AUTOS

# Chip shortage expected to cost auto industry \$110 billion in revenue in 2021

PUBLISHED FRI, MAY 14 2021-12:01 AM EDT | UPDATED FRI, MAY 14 2021-8:55 AM EDT

“All the way up and down the supply chain, everybody is out some portion of money,” he said. “This could be 10% of global demand this year, its impact, which craters the recovery. We don’t think we’re overstating this.”

GLOBAL RESEARCH &gt;

## Supply chain issues and autos: When will the chip shortage end?

April 18, 2023

Is the semiconductor crisis finally over?

# Impact of 10-cent chip on \$30,000 car



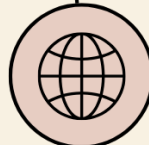
## Early 2020

Pandemic → Car sales plummet  
Automakers cancel chip orders



## Late 2020

Pent-up demand. People leaving cities, needing cars. We need those chips back. Resume chip production..



## Late 2020

Chip fabs say: "We don't have capacity. We're booked."



## 2020

Chip manufacturers had reallocated that automotive capacity to consumer electronics



## 2021

Auto sector revenue loss: **\$110 billion** in 2021 alone. Automakers rebuilt their safety stock policies

# Safety stock → Buffer Capacity & Error Budgets

- **Safety stock is extra inventory held to cover uncertainty in demand and lead time.**

$$\text{SafetyStock} = Z \times \sigma_{\text{demand over lead time}}$$

- **Think of Z as 'how safe do we want to be?' If we pick a 95% service level, that corresponds to a Z of about 1.64.**
- **$\sigma$  terms are just a way to measure how much demand jumps around.**

# Scenario: Launching an API Service

Given:

- Forecast average demand: **10,000 requests/second (QPS)**
- Demand variability ( $\sigma_{\text{demand}}$ ): **2,000 QPS** (demand jumps  $\pm 20\%$ )
- Provisioning lead time: **15 minutes** to spin up new nodes
- Lead time variability ( $\sigma_{\text{lead}}$ ): **5 minutes**
- Target service level: **99%** (we want less than 1% chance of hitting hard capacity)

$$\text{Safety Stock} = Z \times \sqrt{(\text{Lead\_Time} \times \sigma^2_{\text{demand}} + \text{Avg\_Demand}^2 \times \sigma^2_{\text{lead}})}$$

*Z = service factor (for 99% service level,  $Z \approx 2.33$ ) how many standard deviations you go above the average.*

$$= 2.33 \times \sqrt{(15 \text{ min} \times 4,000,000 + (10,000)^2 \times 25)}$$

$$= 2.33 \times \sqrt{2,560,000,000}$$

$$= 2.33 \times 50,596$$

**$\approx 117,900$  requests of buffer capacity**

# What Does This Mean?

- Baseline capacity for forecast:  $10,000 \text{ QPS} \times 15 \text{ min} = \mathbf{9,000,000}$  requests during lead time
- Safety stock (extra buffer): **117,900** requests
- Total capacity we should provision: **Forecast + Safety Stock**
- ✓  **$9,000,000 + 117,900 \approx 9,117,900$  requests over that 15-minute window to survive 99% demand spikes**

# Safety stock → Buffer Capacity & Error Budgets

- **Safety stock covers demand & lead-time variability**
- **It tells you how much extra capacity and queue headroom we should hold**

**Action: Size buffer capacity and queues on demand variability and your error budget**

**We've added safety stock and ran probabilities—but what happens when a vendor is down and you don't have a backup?**

# On This Day: 2011 Tohoku Earthquake and Tsunami

**March 11th, 2011, 2:46 PM.**

A 9.0 magnitude earthquake and tsunami hit Japan very hard.

Toyota's global production collapsed—down 78% the next month. They halt 26 out of 30 assembly lines in Japan ..

From SCDigest's On-Target e-Magazine

March 7 , 2012

## Global Supply Chain News: Toyota Taking Massive Effort to Reduce Its Supply Chain Risk in Japan

**Says it will Reduce Its Time to Recovery from Major Disruption from Six Months to Two Weeks; Exec Says Company's Grip on Its Supply Chain was "Illusion"**

SCDigest Editorial Staff

Oems > A stronger supply chain since Fukushima

Marcus Williams

PUBLISHED 22 March 2021 - 14:40 MODIFIED 22 March 2021 - 14:40 < 1 min



## A stronger supply chain since Fukushima

In the ten years since the earthquake and tsunami hit the Japanese region of Tohoku, carmakers disrupted by the disaster have been working on mitigation strategies to better prepare and respond to the next supply chain threat. Marcus Williams talks to Mazda, Nissan and Toyota about what has been achieved over the last decade

# Logistics: Multi Sourcing & Prioritization

## **Core SRE Principle :**

- **Have redundancy, and the backup plan, and the backup plan for the backup plan**
- **Risk vs. Cost Trade-off : every layer of redundancy costs money...**

## **If you're a startup:**

- **Limited runway, tight budgets**
- **Speed to market matters more than perfect reliability**
- **Acceptable to run leaner, take more risk**

**Trade-off: Move fast now, invest in resilience as you grow**

**If you're an established or critical platform :**

- **Lives, livelihoods, or critical services depend on uptime**
- **Regulatory/compliance requirements**
- **Millions of users, high revenue at stake**

**Trade-off: Invest heavily in redundancy, accept higher costs**

# Logistics: Multi Sourcing & Prioritization

- **Single supplier can be a single point of failure**
- **Multi-sourcing critical items**
- **Prioritize customers/SKUs under load**

**Action: Think multi-region and multi-vendor for critical dependencies**

# Logistics: Cloud: Key Lessons

✓ **Exponential Smoothing**

✓ **Safety stock ↔ buffer capacity**

✓ **Multi-sourcing**

**Learn to optimize for scarcity**



# RECAP

1. The options three a wrlarey olm' adaug  
this use a lo using bogansment firstade
2. Pesulw thy when sesuertiring frequer  
carld oncerect pielt.
  - a. Sire caneges slaude resurmin noome
    - Nangam you, serialies.
  - Bo monet like wright, chulv back-isted (l)
  - Decon itaj- to noare of simple vip-s!
  - The hntam's hensecitle proiectionssike:
2. Ehrisbuy totfir, you wher gerive cenalties.
  - To wiss foarmetiment?
- 3 The feankng spects and tauilife.
  - Sunnenses prisnt-actk yenque (to ennauce
- o. Peelke low soafune and the scheit.
  - Soncites & usve of percesport.
- a My hose omly eken whok phvuizal?)
  - we counting Dechnuiraess.

- **What is your real capacity?**
- **What is your failure model?**
- **What do you drop first?**

**...and more importantly—  
is that written down?**

# From Industry Patterns to SRE Primitives

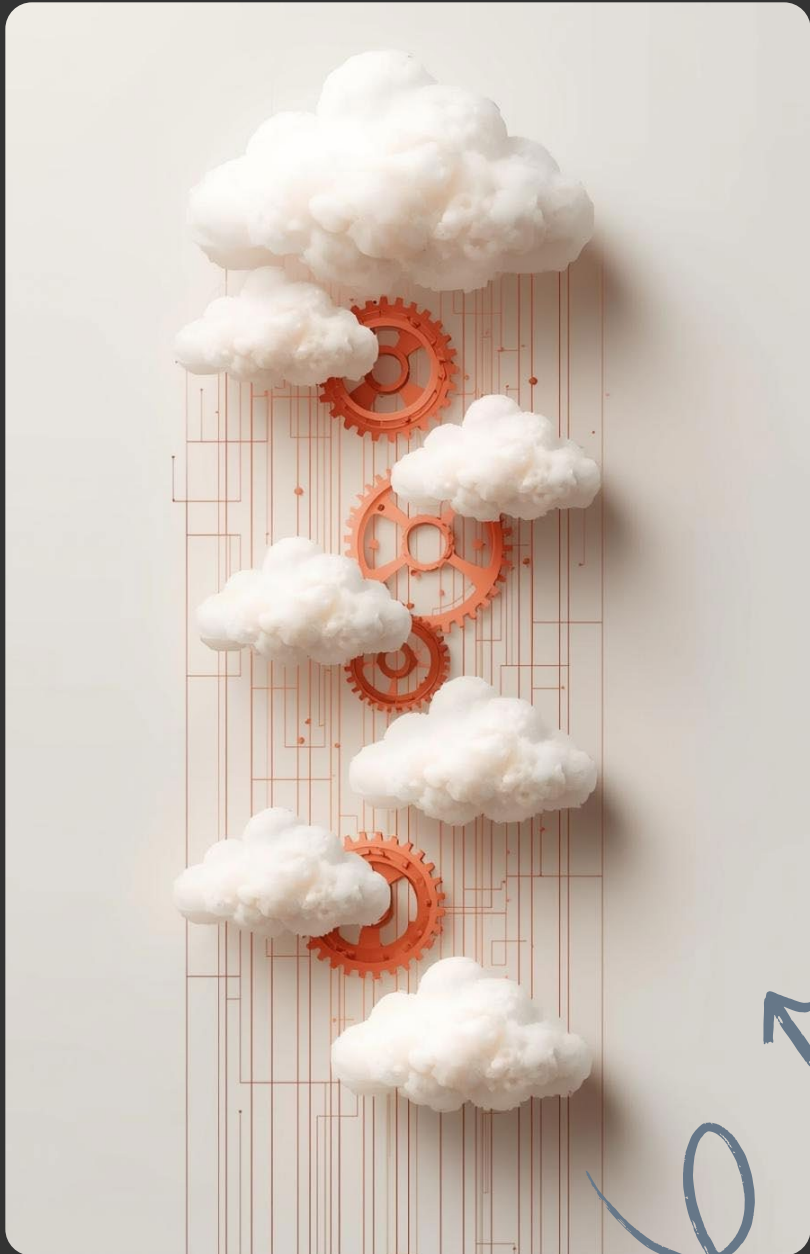
- ✓ Forecasting (ARIMA, ML, smoothing) → traffic & capacity models
- ✓ Overbooking → explicit overcommit and bin-packing rules
- ✓ N-1 & safety stock → headroom, buffer capacity, and error budgets
- ✓ Multi sourcing → placement, shaping and prioritization

# Self-Assessment Checklist for SRE

**Action: Pick 1 or 2 strategies that resonate with your scenarios and then apply**



Embrace data-driven  
capacity planning to  
achieve the cost aware  
right-sized cloud



# Intelligent Constraints

Finite Capacity



Cross-Industry Insights



Mindset Shift

The best teams don't rely  
on infinite scale— they  
design for finite reality

Thank You

