



Building warehouse-scale computers

or ... what's it like to supply exponential growth

john wilkes 2018-10



~~You're all~~
Some of you are
thinking
too small

Scale has been the single most important force driving changes in system software over the last decade

– *Technical perspective: Is scale your enemy, or is scale your friend?*

John Ousterhout, CACM 54(7):110, July 2011.

The kind of things we like to be able to do

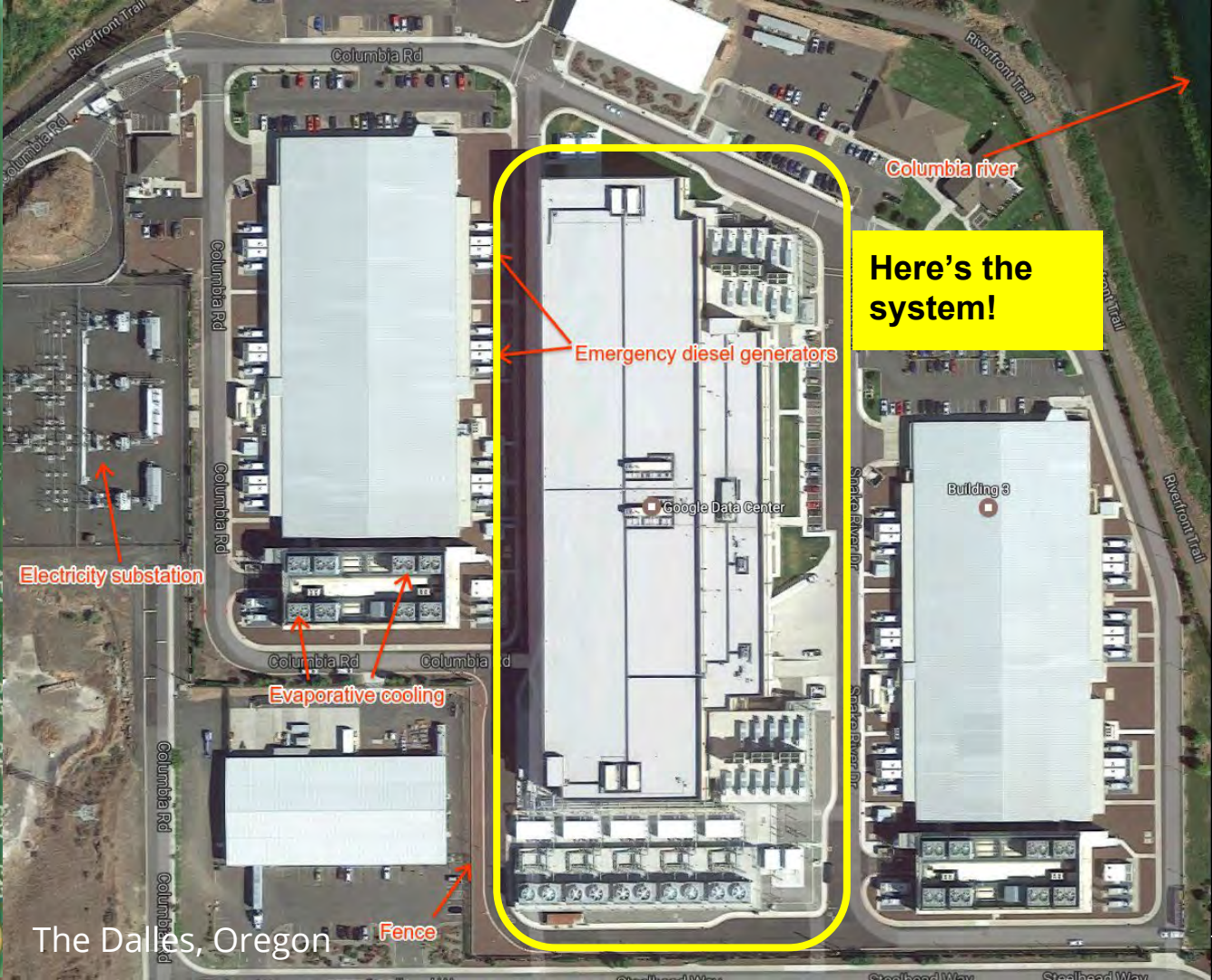
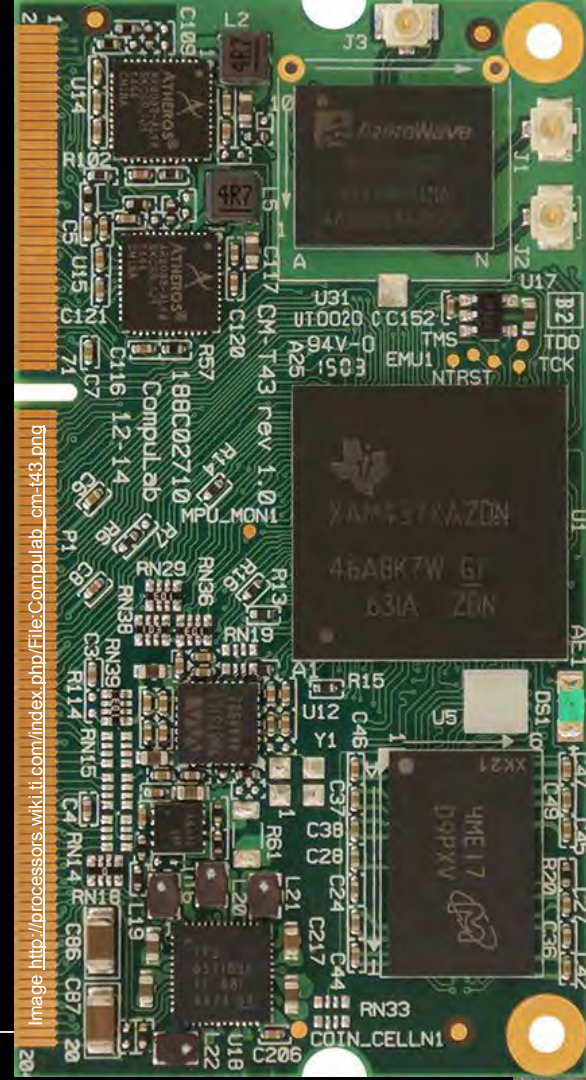
The screenshot displays the Google BigQuery web interface. On the left, the 'COMPOSE QUERY' sidebar shows a tree view of datasets, including 'google.com:bigquery-petabyte' and its sub-dataset 'retail_petabyte'. The main area shows a 'New Query' editor with a SQL query. Below the query editor, a red 'RUN QUERY' button is visible. To the right of the query editor, a pink box highlights a status message: '1PB table; 13.3s elapsed; 266GB processed'. Below this, a red box highlights the message 'Query complete (13.3s elapsed, 266 GB processed)' with a green checkmark. The bottom of the interface shows a table of results with columns 'Row', 'sale', and 'day'.

```
1 SELECT INTEGER(SUM(totalSale)) as sale, DATE(orderDate) day
2 FROM [google.com:bigquery-petabyte:retail_petabyte.sales_partitioned]
3 WHERE customerKey = "1104796800000-155"
4 AND PARTITIONTIME BETWEEN
5     TIMESTAMP("2005-01-04") and
6     TIMESTAMP("2005-02-04")
7 GROUP BY day
8 ORDER BY day desc
9
```

1PB table; 13.3s elapsed; 266GB processed

Query complete (13.3s elapsed, 266 GB processed)

Row	sale	day
1	622335699	2005-02-04
2	623106720	2005-02-03



The Dalles, Oregon

Here's the system!

Emergency diesel generators

Evaporative cooling

Fence

Electricity substation

Columbia river

Google Data Center

Building 3

What if ...

the “**system**” included not just the computers,
but the network fabric and the WAN endpoints



and the cooling system



and the building-management system



St Ghislain, Belgium

and the power system
and the technicians



Backup generator at
St. Ghislain, Belgium

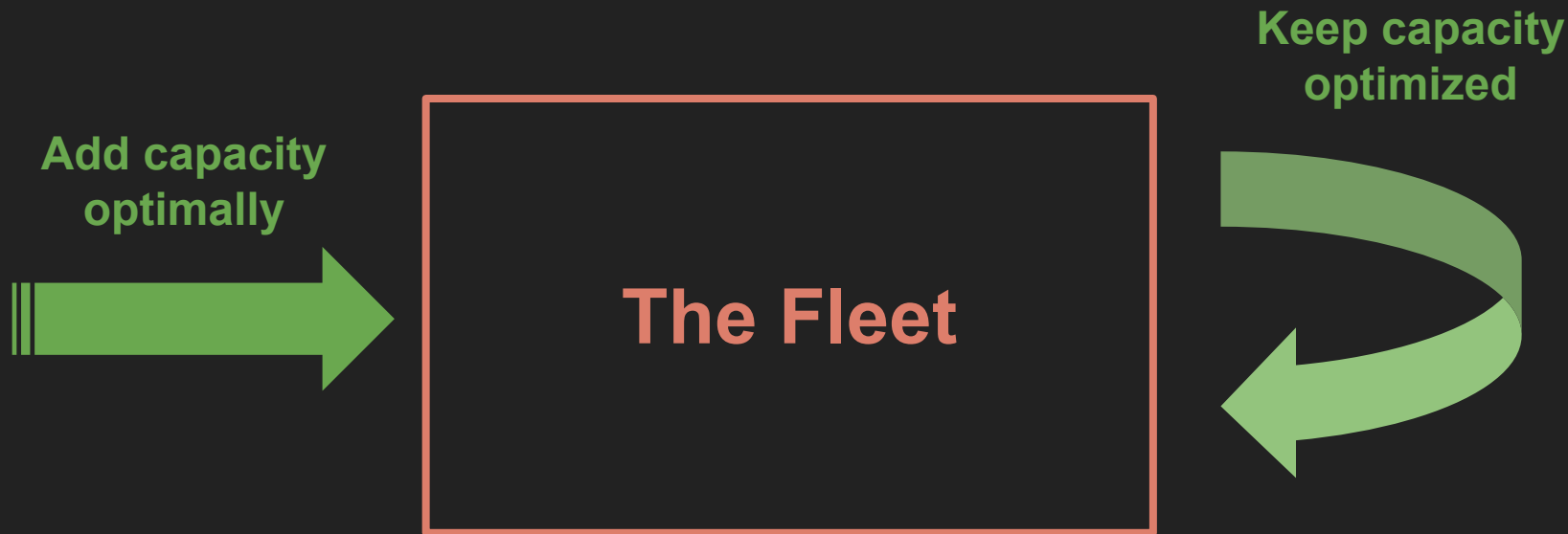
In fact ... the entire **warehouse scale computer**



By the way ... it costs O(\$200M). No pressure.



Key challenges



Planning for compute resources

All have multiple generations/options

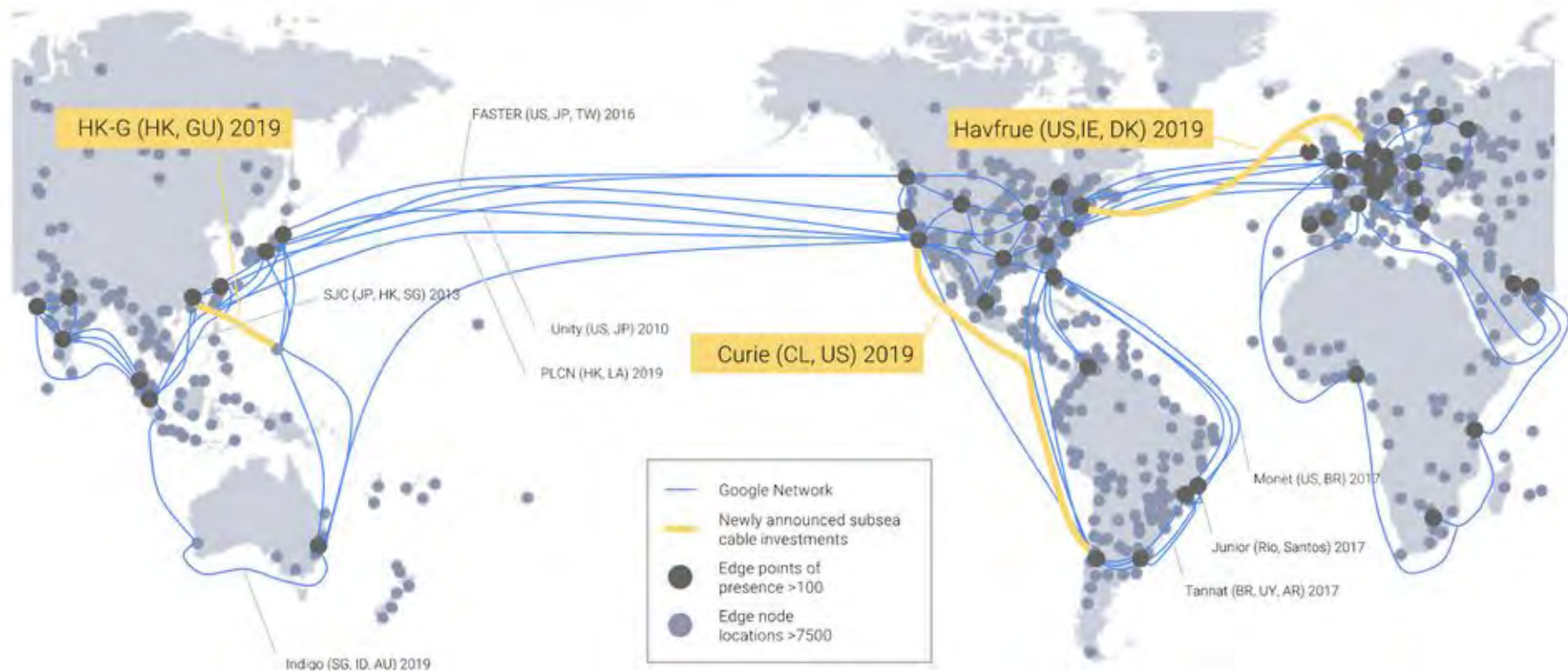
- CPU
- RAM
- disk, SSD – capacity, performance
- new NVRAM thingies ...
- accelerators (complication: a wide variety)
- inter- and intra-datacenter networking
- power
- datacenters
- land, water, sewage
- ...

Response times are weeks/months/years, not milliseconds



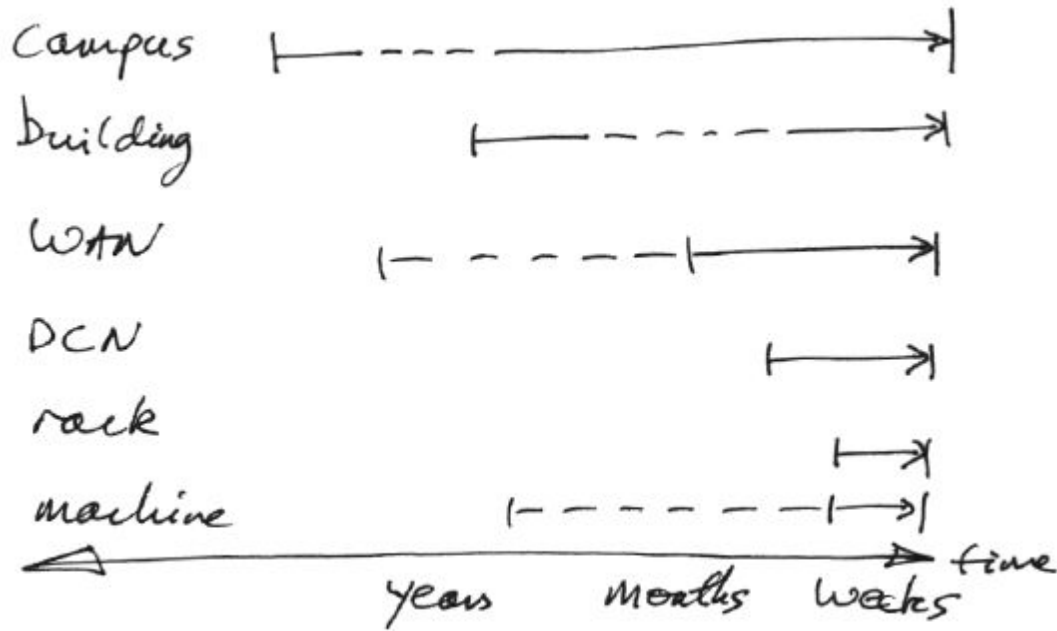
Google Network

The largest cloud network, comprised of more than 100 points of presence



How long does it take to build capacity?

Note: variance is as bad as delay



Planning for compute resources

What kind of machine to buy?

- different groups want different things (e.g., search, Cloud)
- ⇒ MotD + customizations

Idea: reuse machines when they get handed down

- when is it worthwhile? (price of power? room for expansion?)

Also: do we have ... space? power? networking? budget?

- models; what-if analyses; uncertainty (demand, supply)
- objective function: *total cost-of-ownership*

Planning for compute resources: Total Cost of Ownership (TCO)

what do you mean by “**total**”?

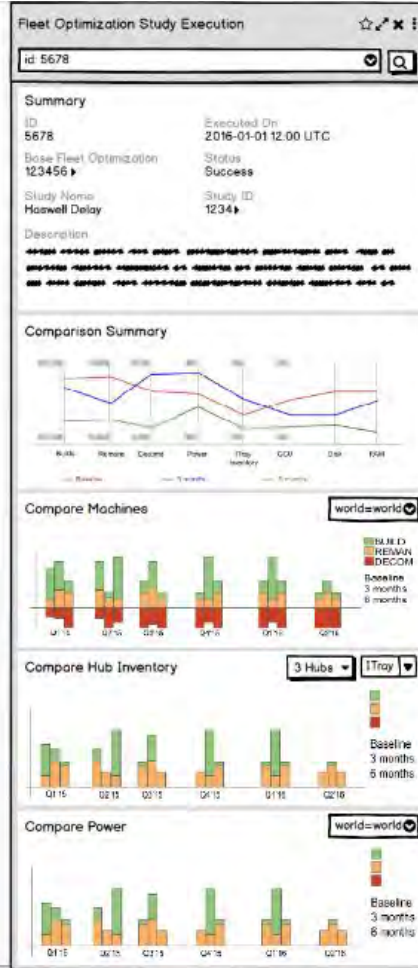
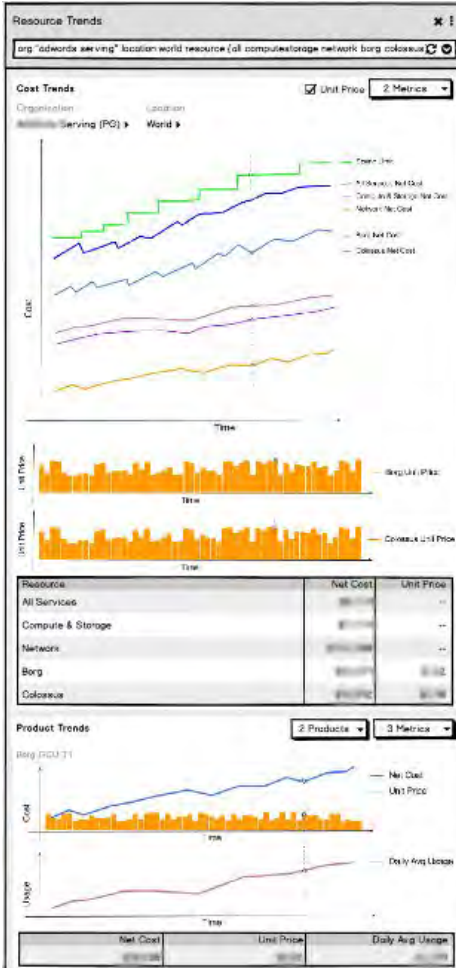
- average? over what? continent? time?
- depreciation schedule?

what do you mean by “**cost**”?

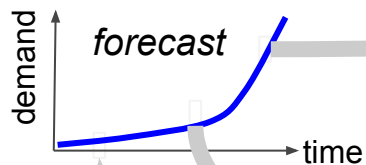
- initial purchase price, or the average over time?
- is delivery or installation included? what if they are being reused?

what do you mean by “**ownership**”?

- machines can be transferred/given/sold/break
- who “owns” a machine running a shared service for a customer?



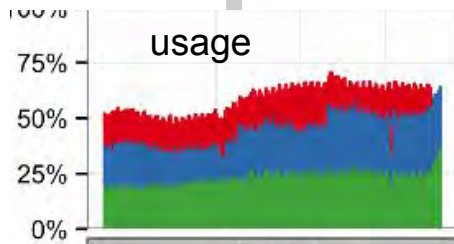
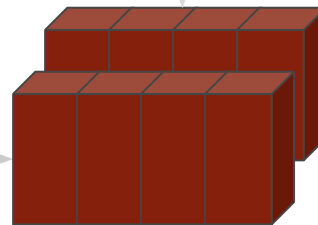
A simplified overview



sites, data
centers, power



network,
machines,
storage

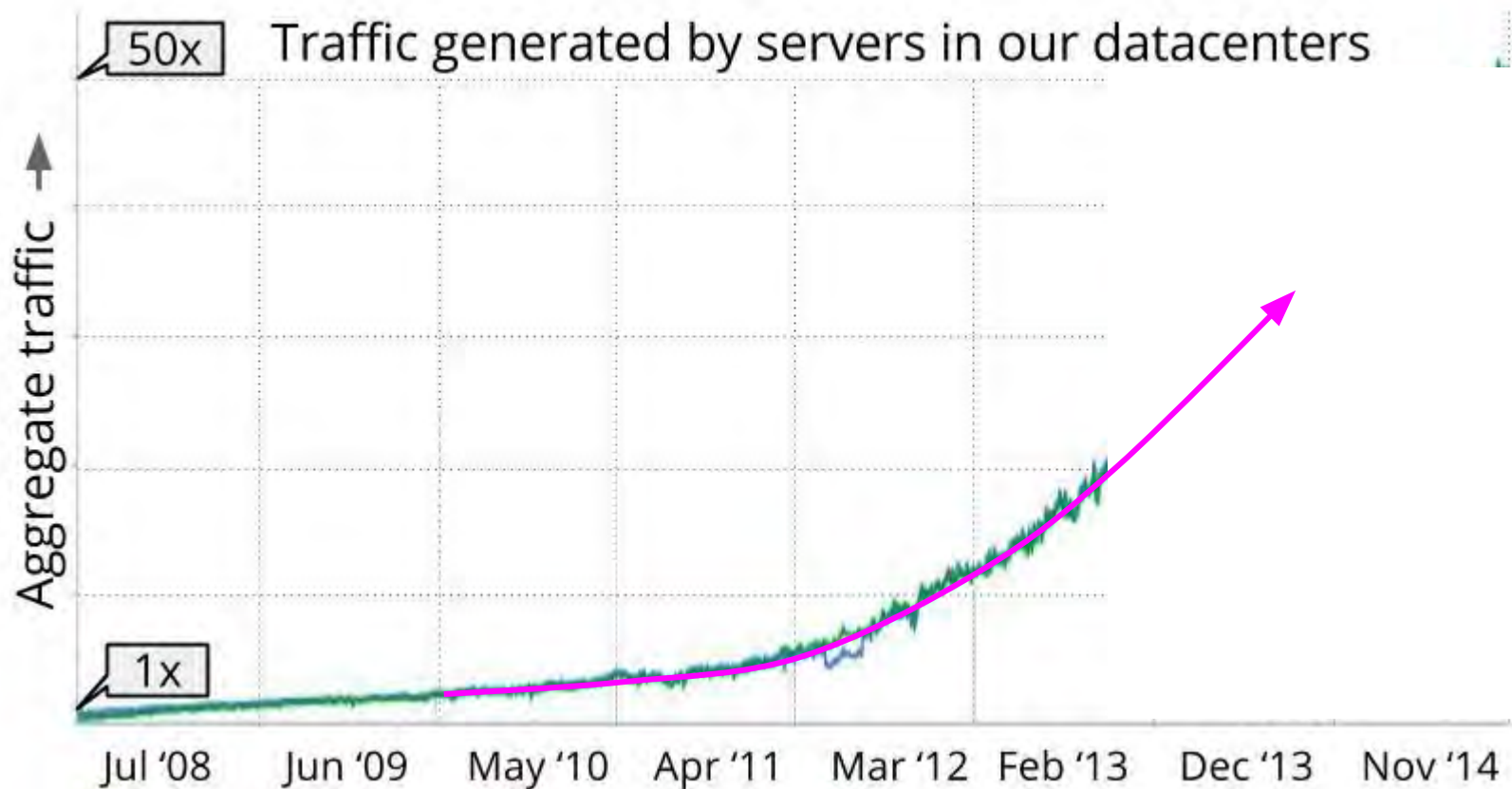


\$\$ prices

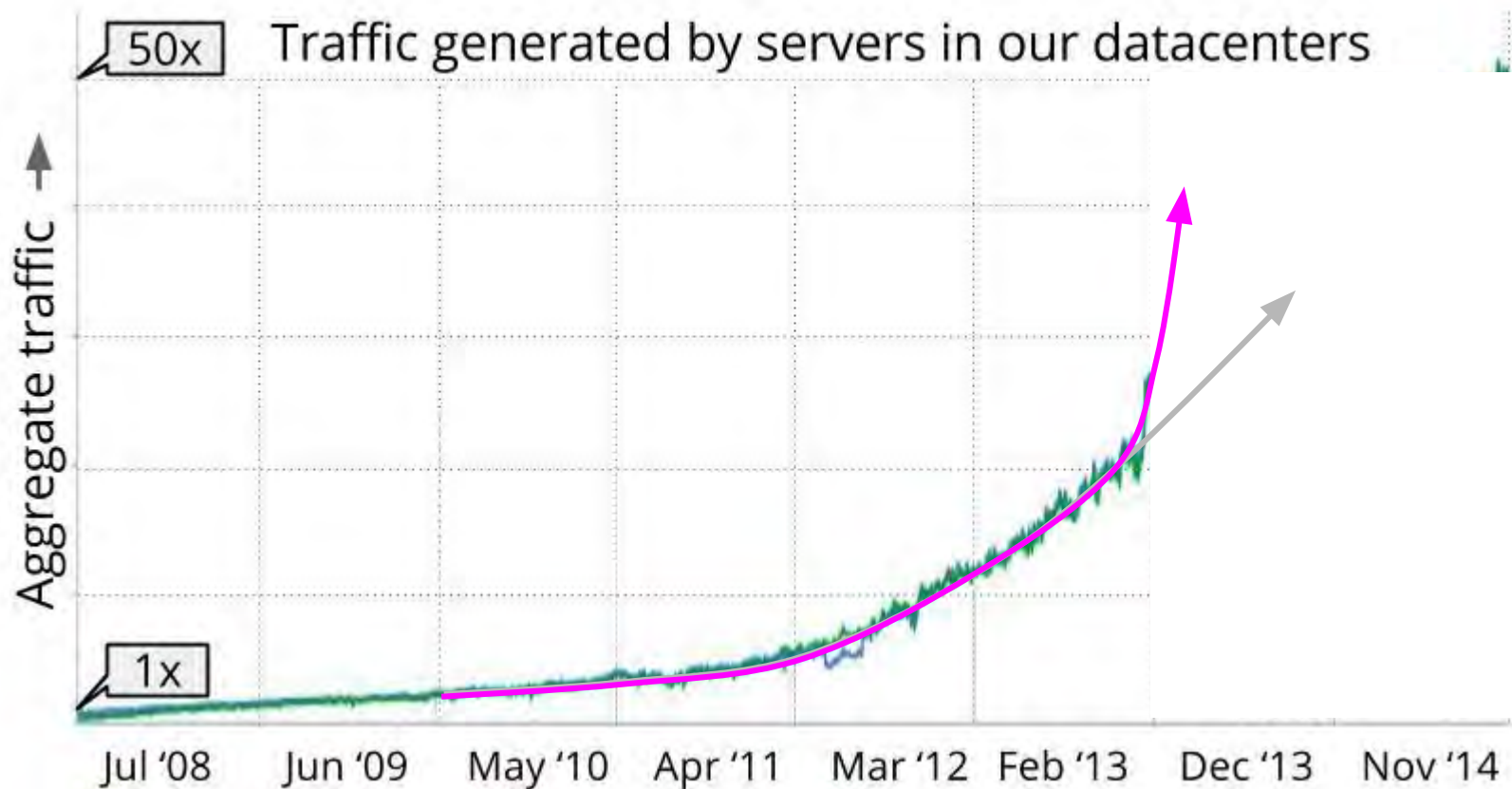


compute/storage
capacity

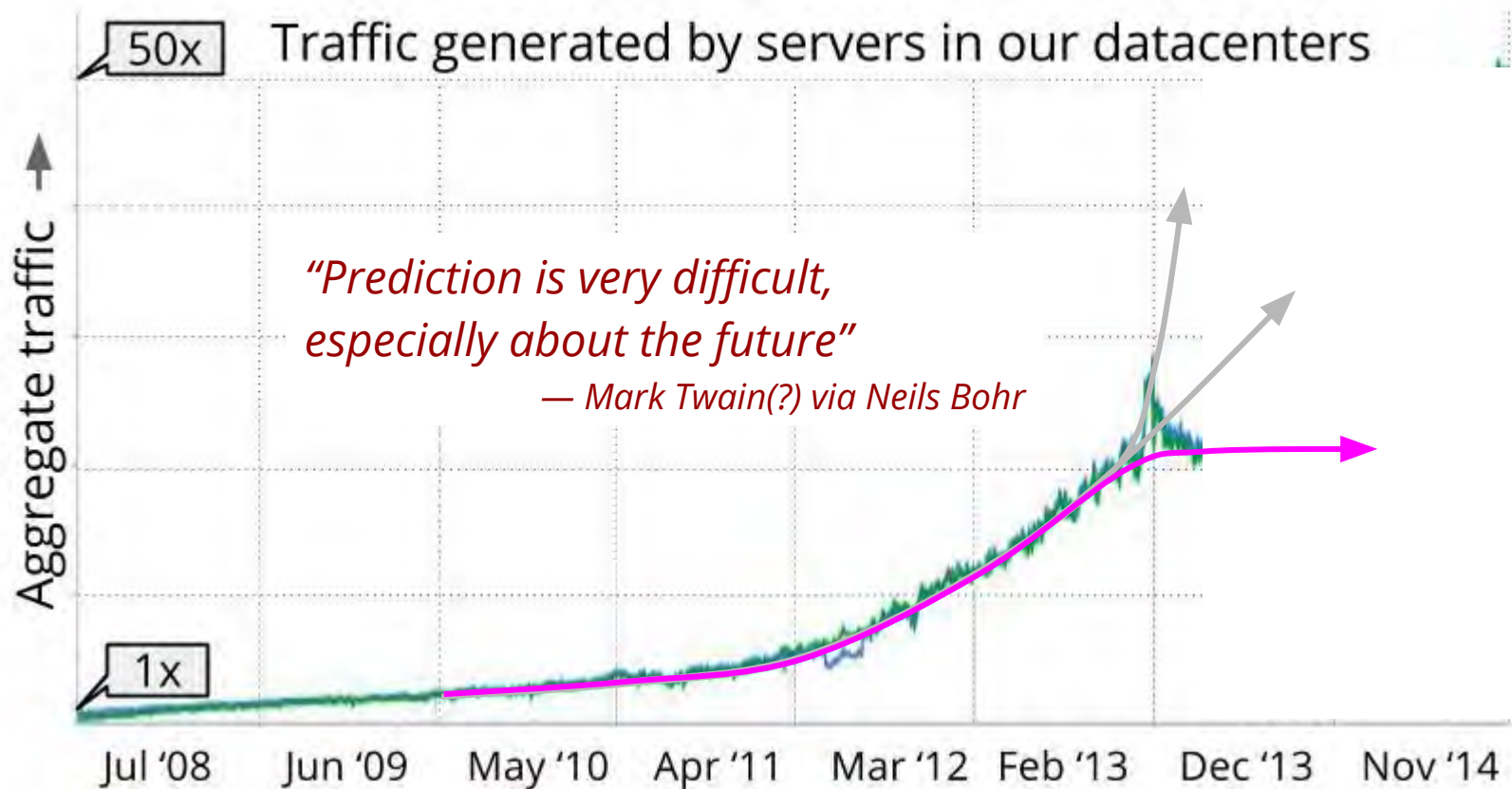
In a world of exponential demand growth



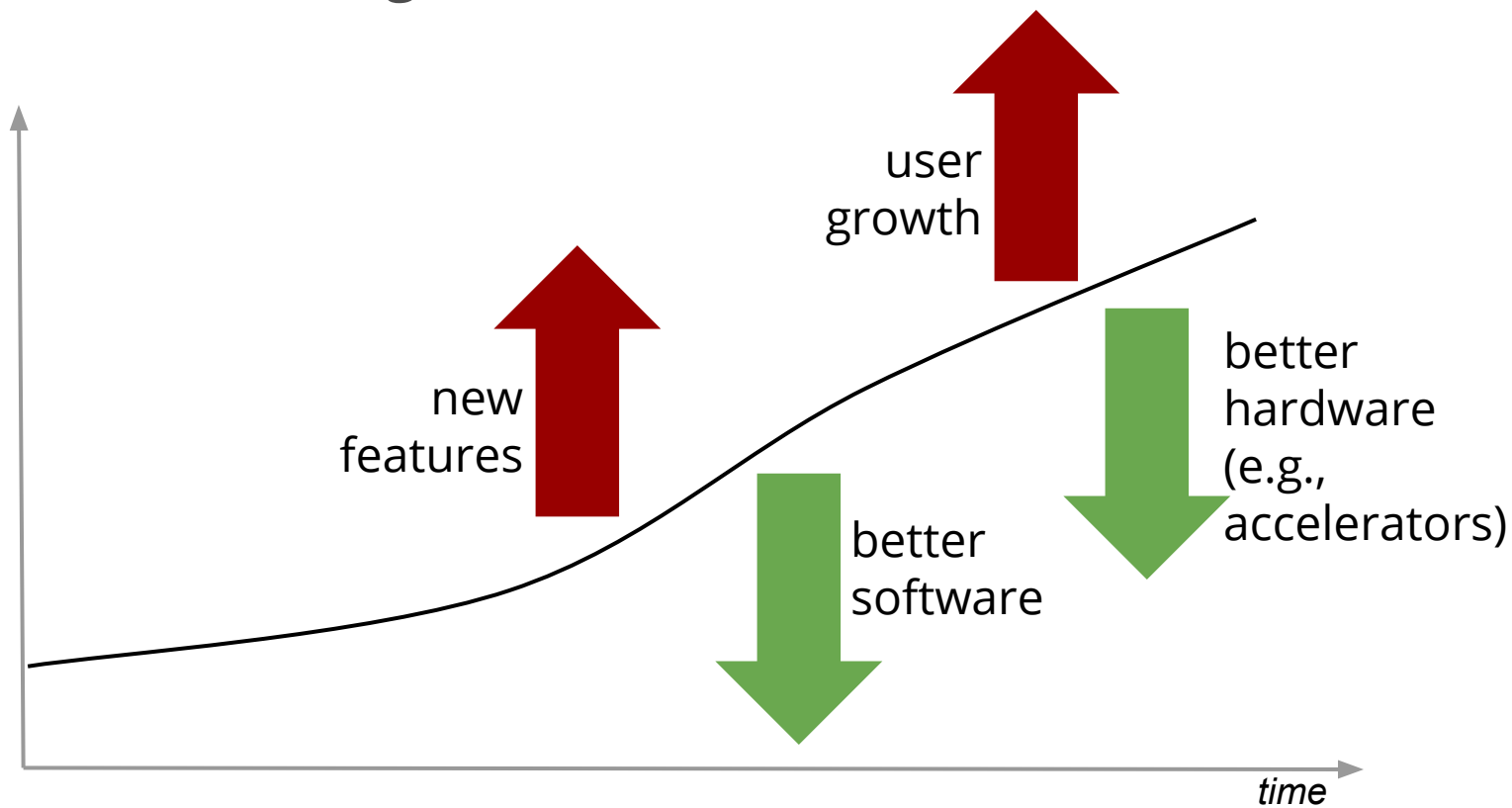
In a world of exponential demand growth



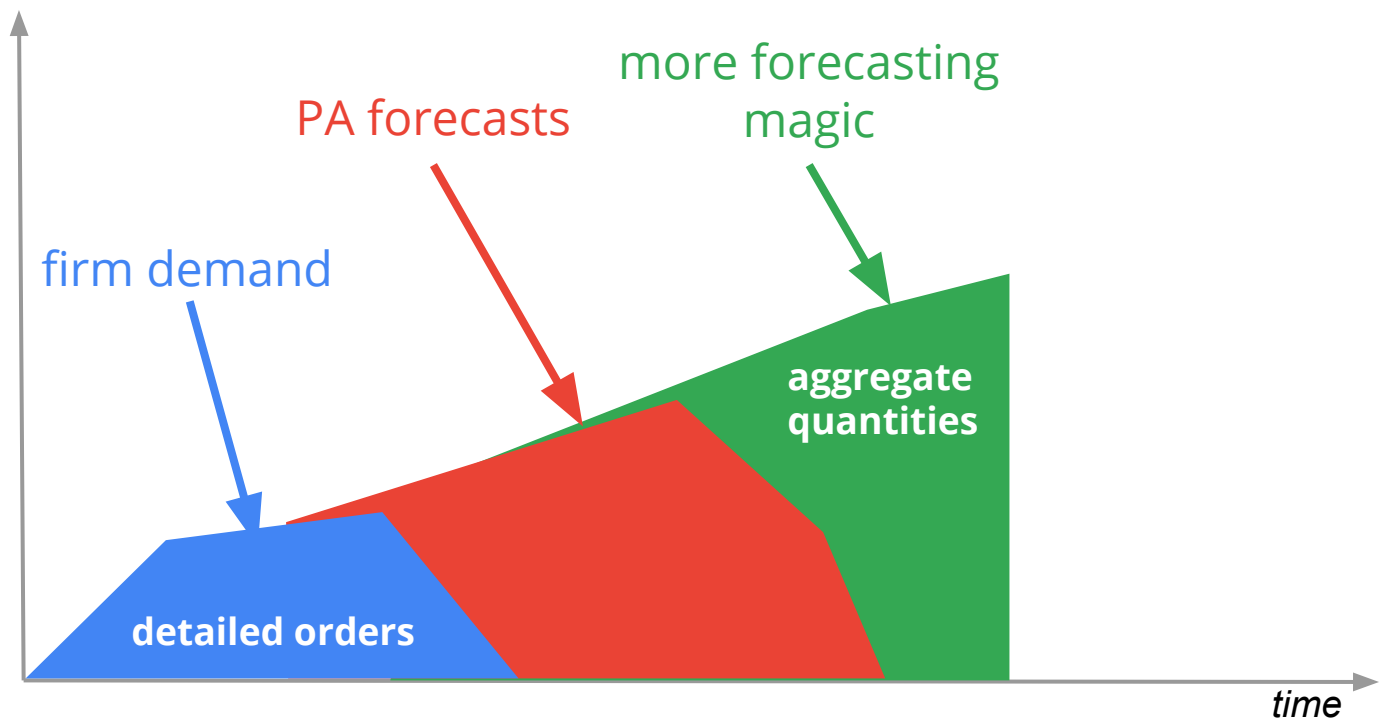
In a world of exponential demand growth



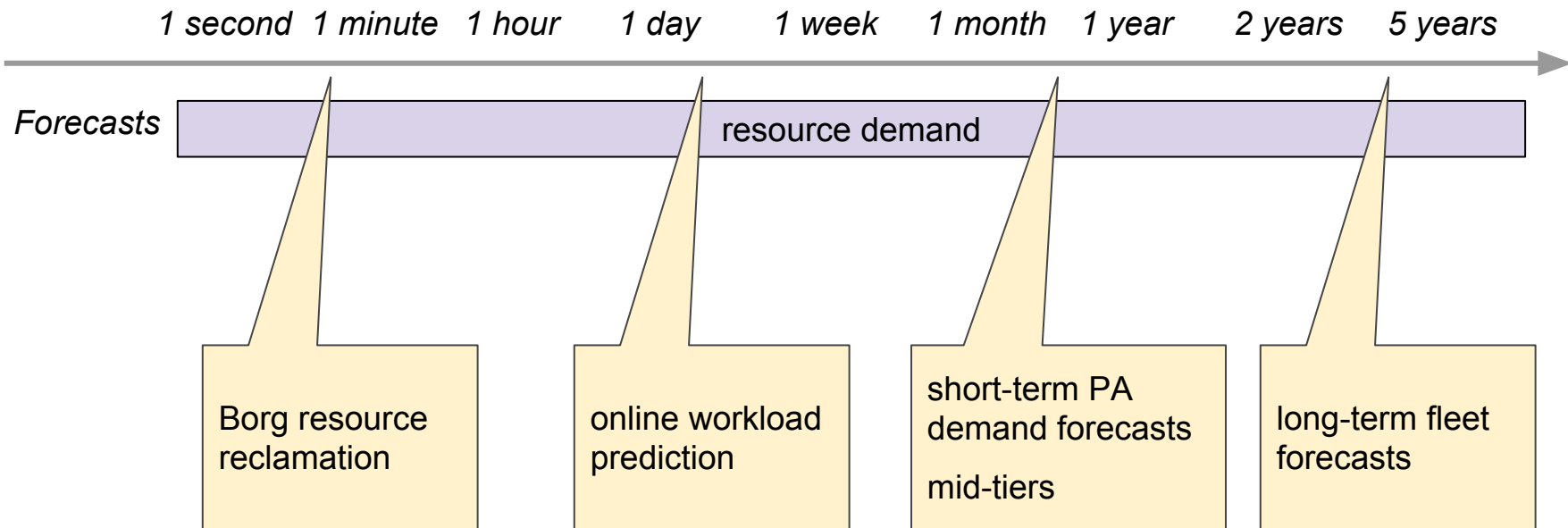
A few factors affecting forecasts



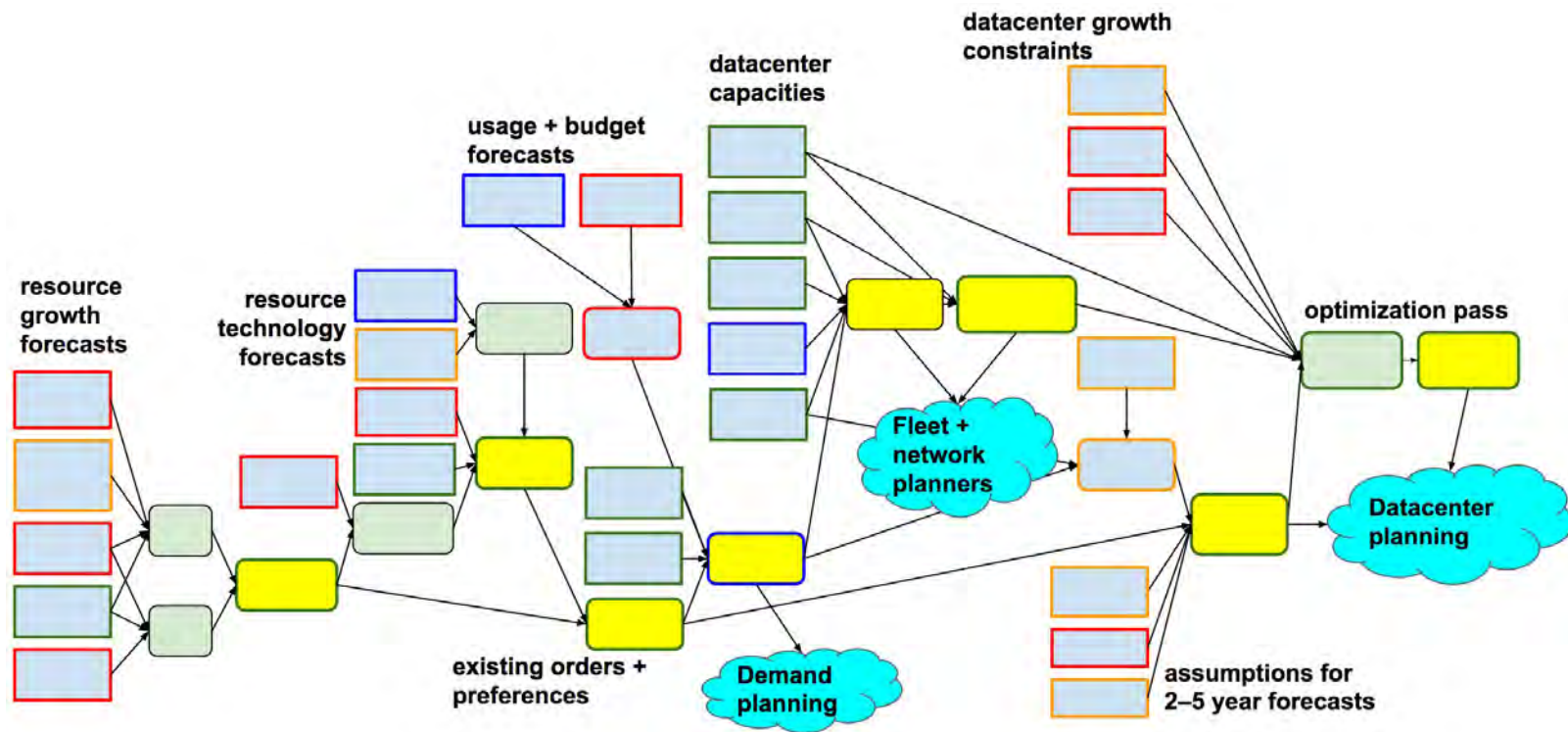
How much capacity do we need? *and when do you need to know?*



Putting it all together



It takes a few moving parts ...



1X
Target Traffic

5X
Worst Case
Estimate

50X
Actual Traffic

Cloud Datastore
transactions/s



Planning for power

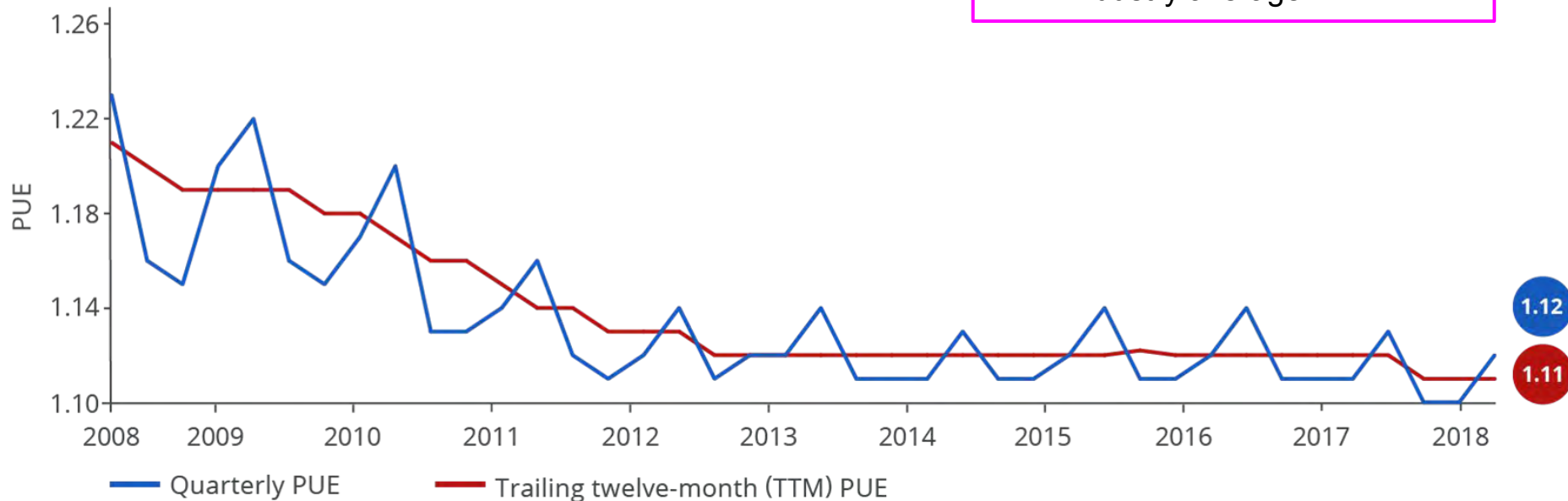
Continuous PUE Improvement

Average PUE for all data centers

PUE (Power Usage Effectiveness)

= *total Watts / compute Watts*

- smaller is better
- industry average ~1.7



Planning for power



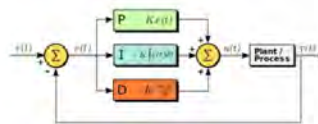
Planning for power

Some of the challenges

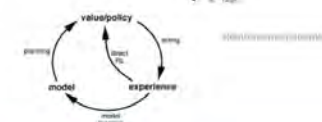
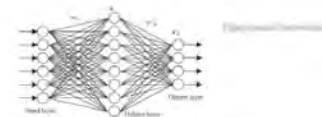
- Narrow range of experience
- Exploration may discover unsafe states
- Inputs out of our control (e.g., weather)
- Control system reliability/availability
- Agility (new hardware)
- Reinforcement learning with long delays between action and change in system state
- **Safety. Safety. Safety.**

Tier-2 control systems

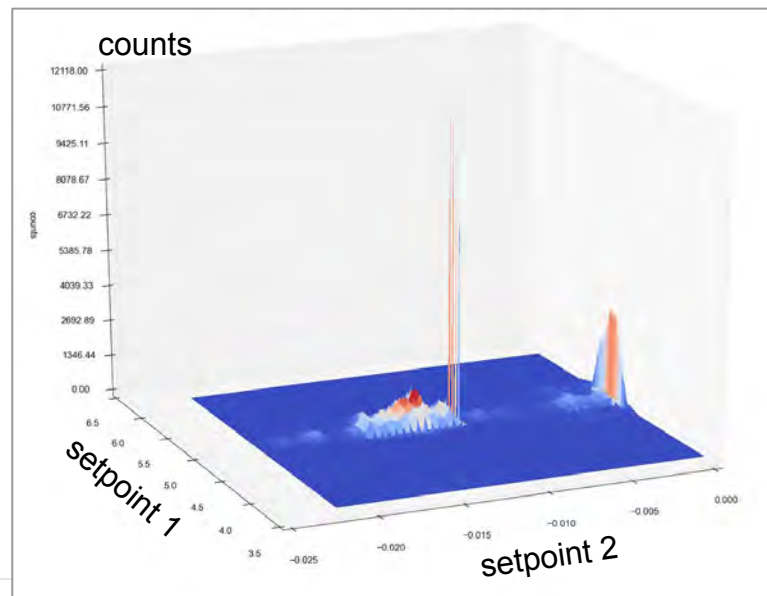
- Layered on Tier-1 critical systems
- Typically provide efficiency
- Designed to fail safely to Tier-1



Classical control in the cloud



ML-based control in the cloud



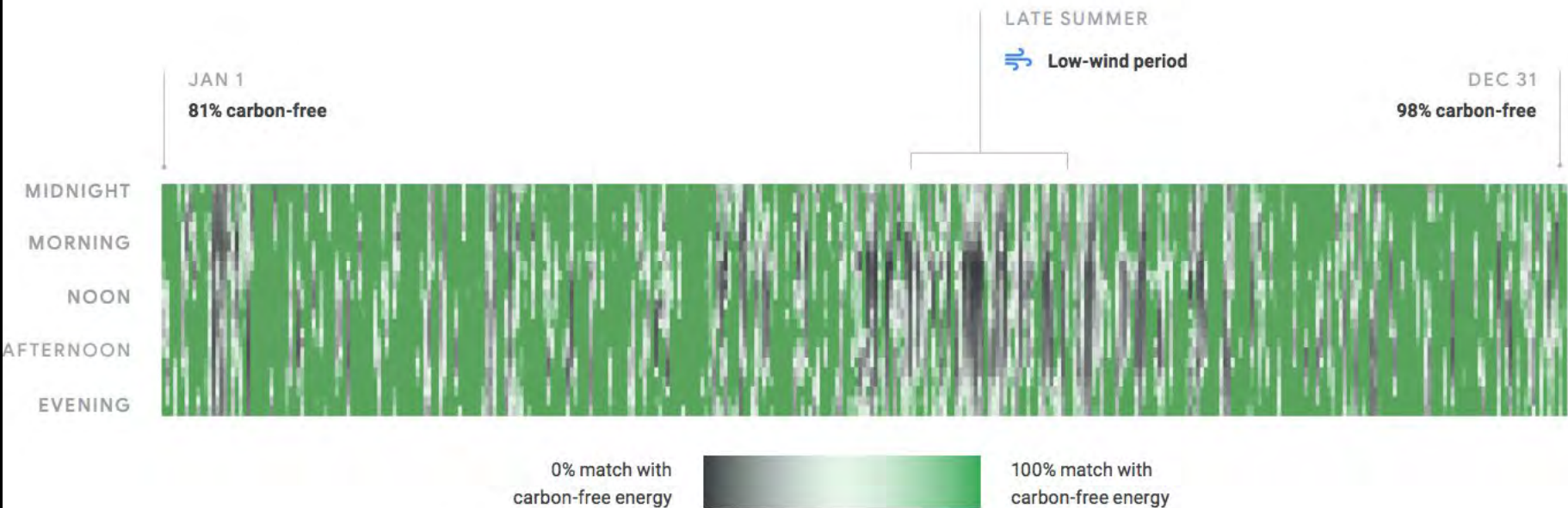


We're set to reach 100% renewable energy in 2017

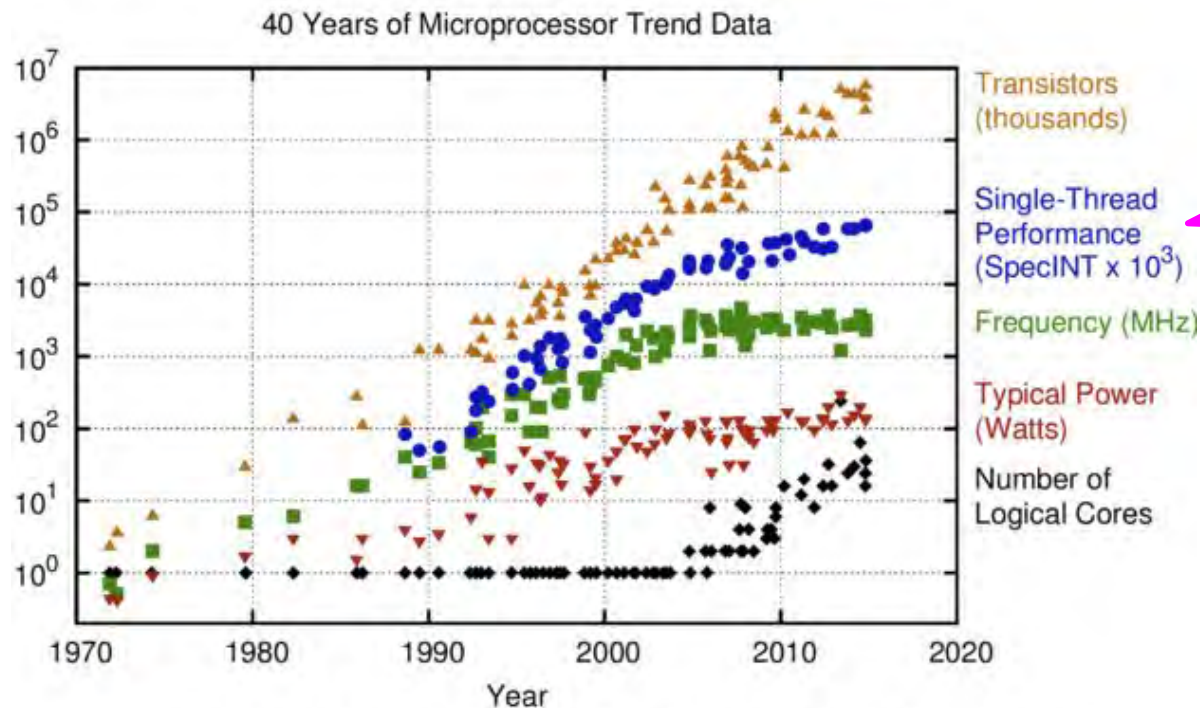
Every hour of electricity use at Iowa data center

Although our Iowa data center achieved 100% carbon-free energy during the majority of hours in 2017, there is also a recurring reliance on carbon-based power — most notably in late summer, when wind speeds decline.

Overall in 2017, 74% of this data center's electricity use was matched on an hourly basis with carbon-free sources.



Meanwhile – what's up with Moore's law?



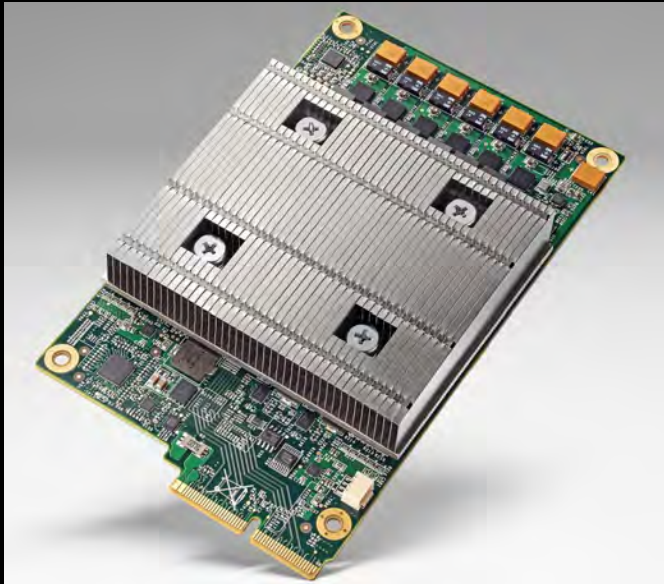
Original data up to the year 2010 collected and plotted by M. Horowitz, F. Labonte, O. Shacham, K. Olukotun, L. Hammond, and C. Batten.
New plot and data collected for 2010-2015 by K. Rupp

Single-core performance plateauing after decades of exponential growth

Graph from [40 Years of Microprocessor Trend Data](#), Karl Rupp, CC-BY 4.0.

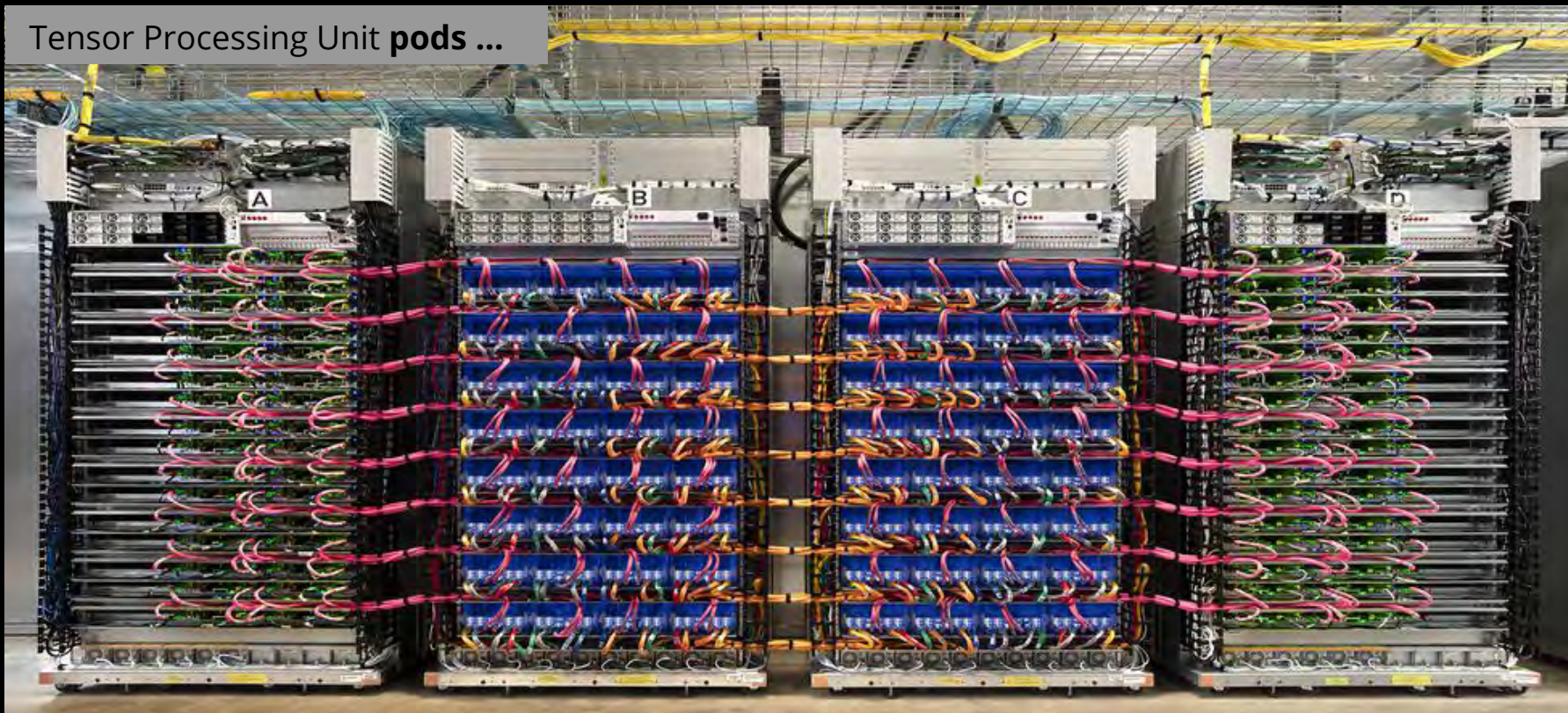
Meanwhile – what's up with Moore's law?

Tensor Processing Units

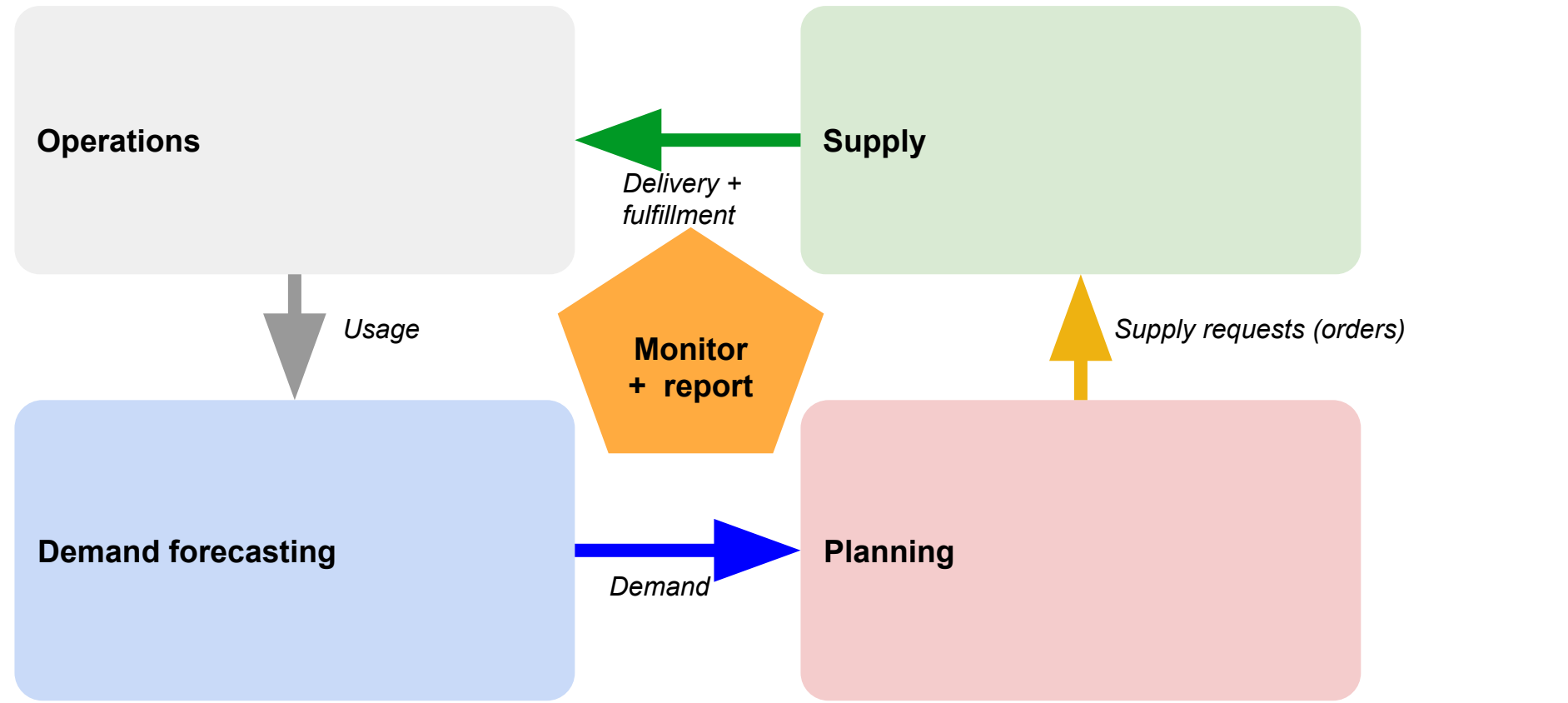


Meanwhile – what's up with Moore's law?

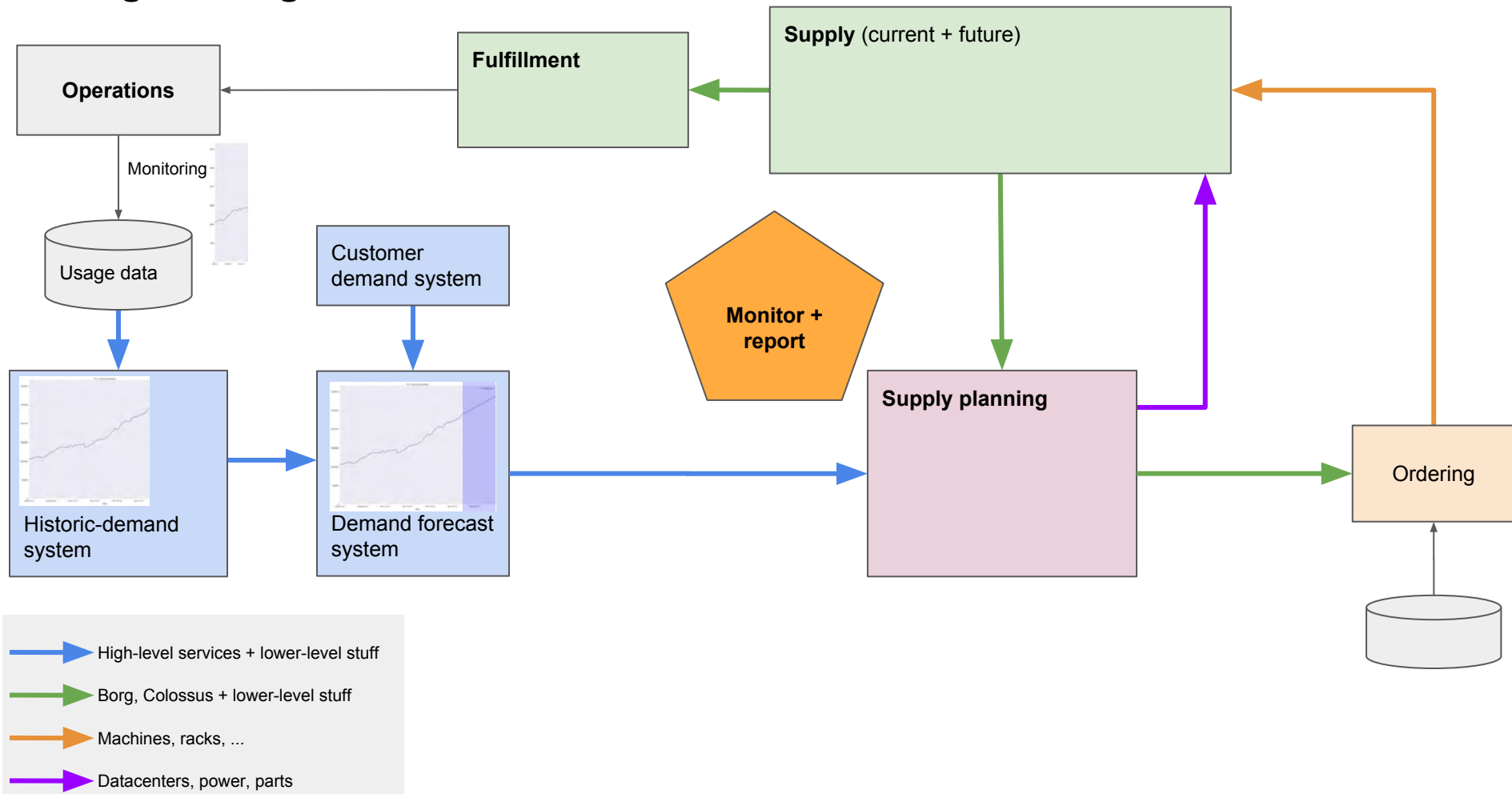
Tensor Processing Unit **pods** ...



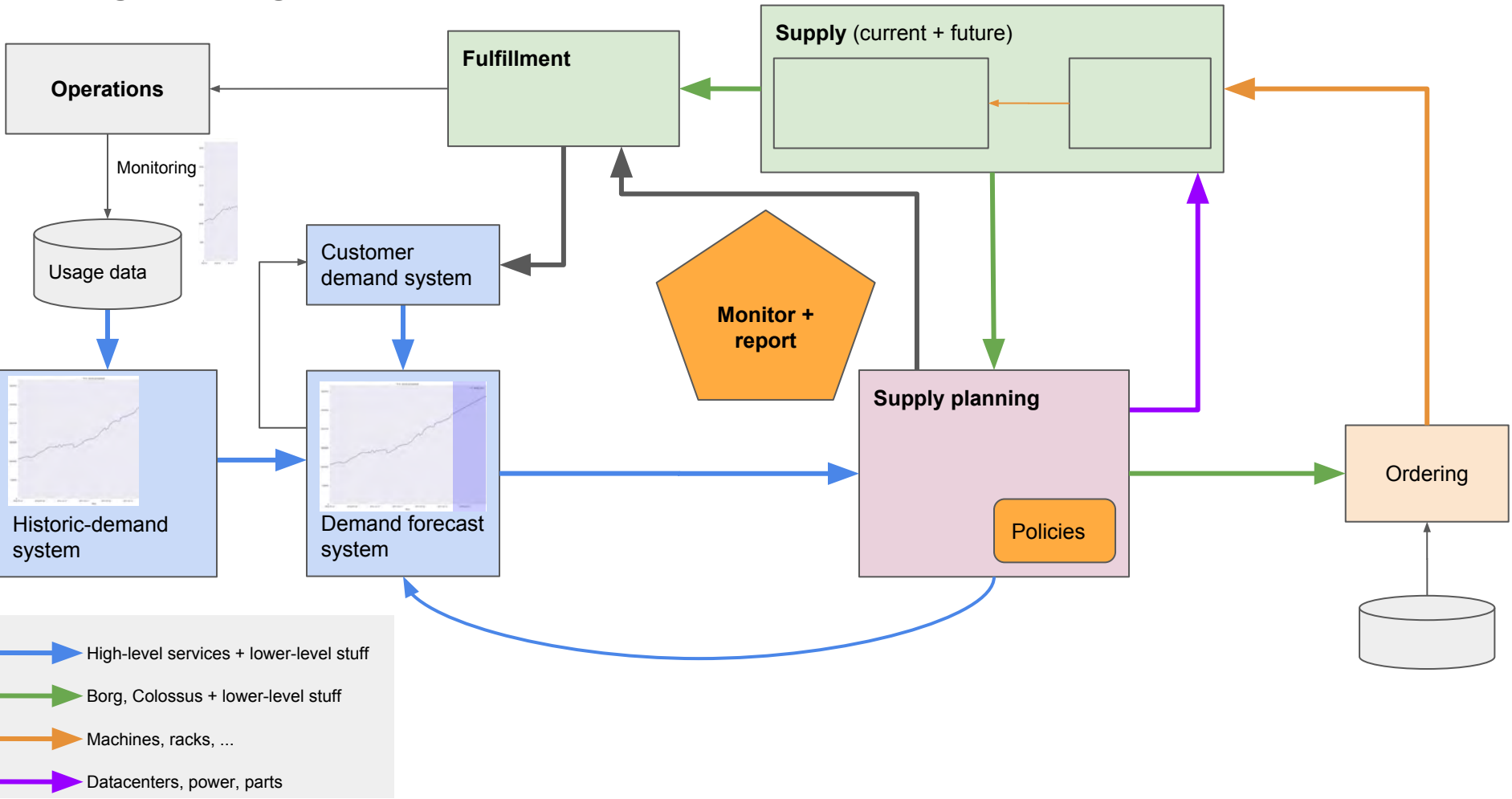
Putting it all together



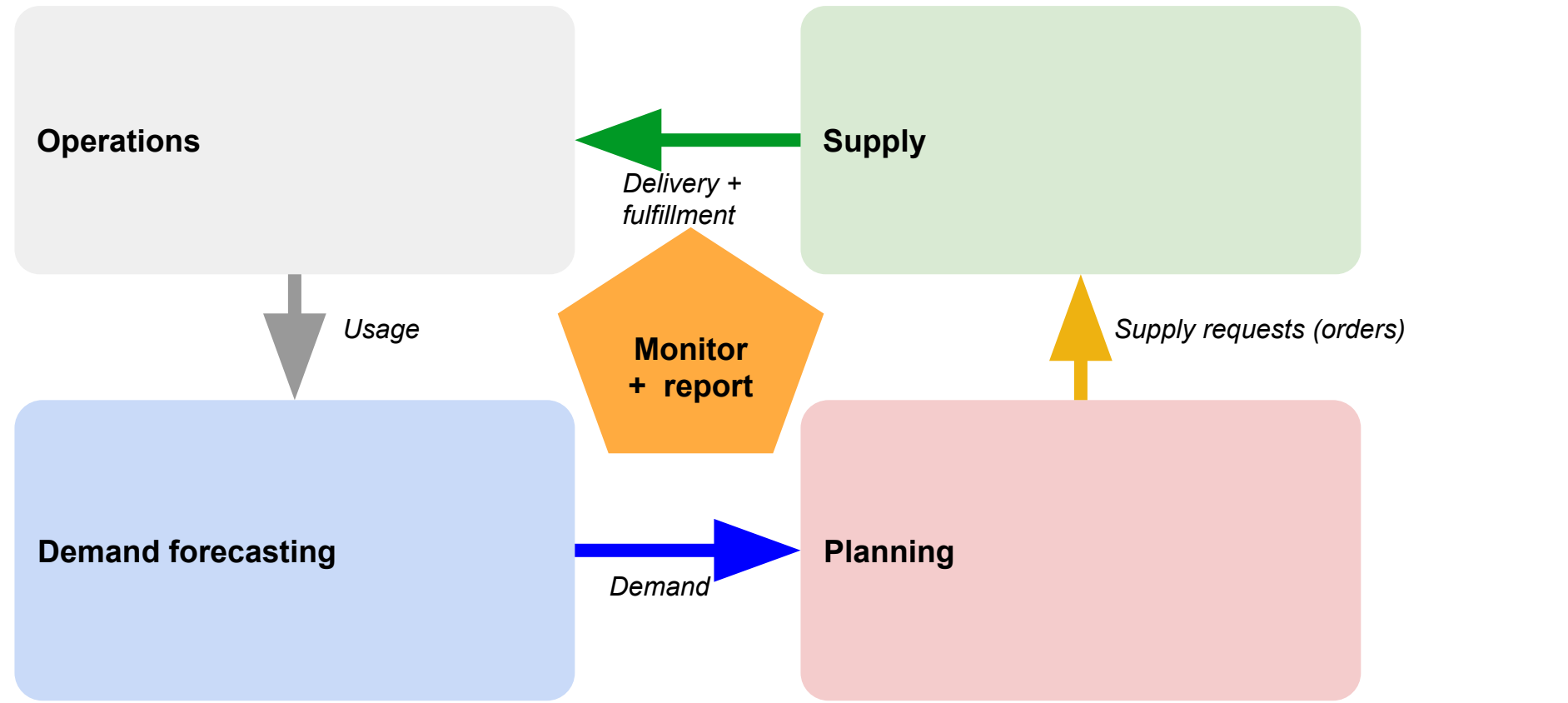
Putting it all together – a few more details



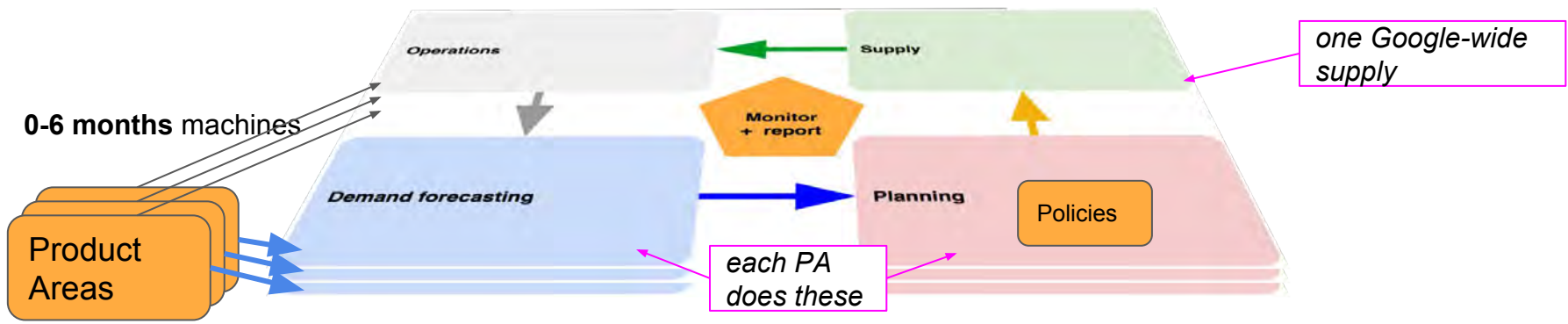
Putting it all together – a few more details



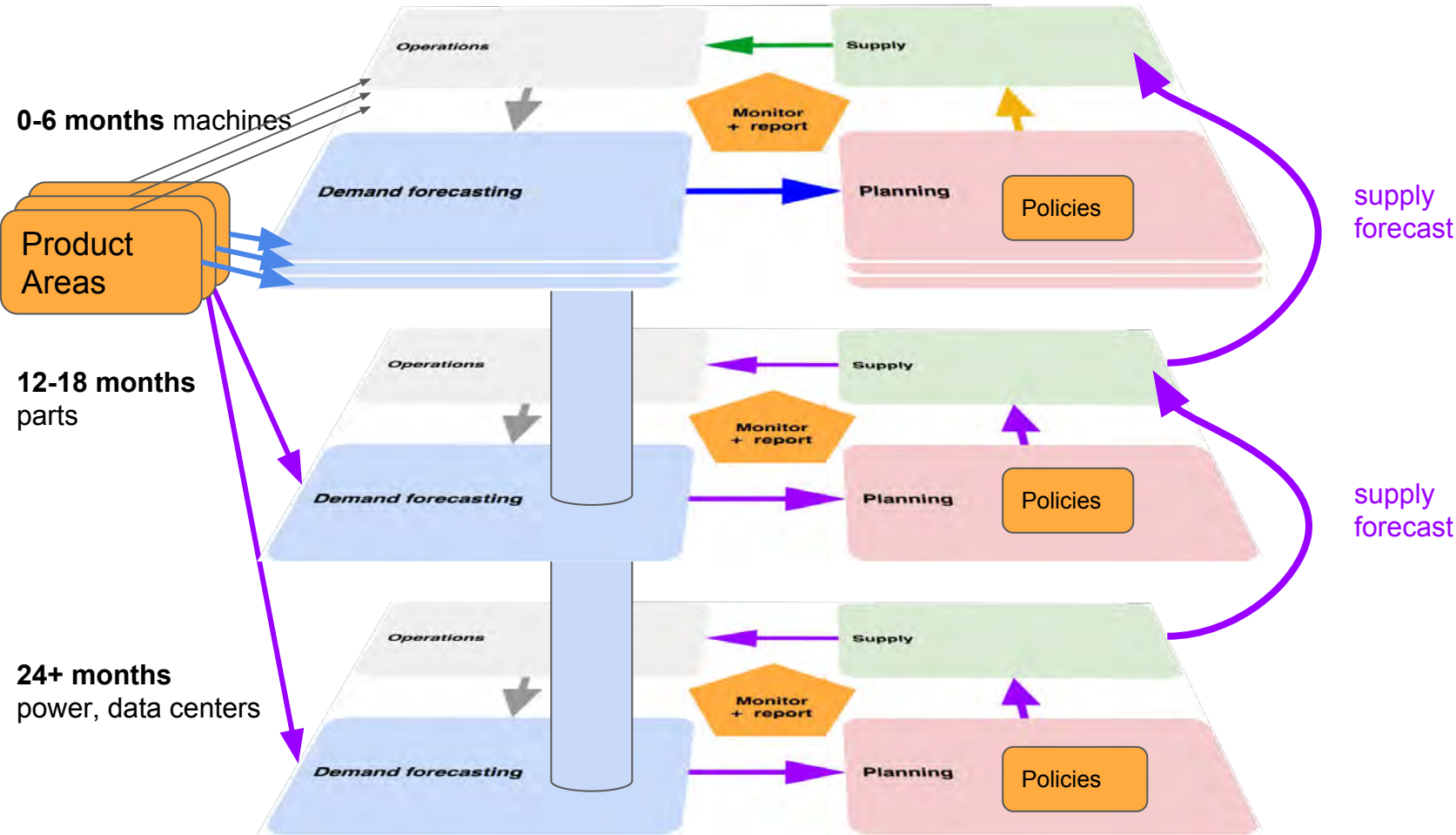
Putting it all together



Putting it all together



Putting it all together – multiple timelines





2018 Q1 CapEx = **\$5.3B**

(+\$2.4B for an office building in New York)

source: Alphabet [SEC filing](#)

\$29.4B

3-year trailing CapEx, as of March 2017

Final thoughts:

There's a lot of technology
behind "the cloud"

At scale, efficiency *really*
matters

[and: we're hiring!]

johnwilkes@google.com