**Microsoft**

# *Optics for Cloud:* new approaches to data centre technology

EPSRC CDT in Applied Photonics Annual Conference, Heriot-Watt University, 7 June 2019

Tom Empson
Azure Storage Research Laboratory

*Dublin*

# Data is doubling every two years

**2013**

**2020**

4.4 ZB

**44 ZB**

Slide credit: Mark Russinovich
Source: IDC 2014

# 1 ZB is one billion TB: we're facing a data tsunami

- Current HDDs can store 10 TB...

- So we need 100 million of those to store 1 ZB...

- That means 1000 data centres...

- Which would cover 20% of Manhattan

- This data also needs to be transmitted across the network...

- Data centre network links can send 100 Gbps...

- That's 2500 years to transmit a zettabyte

**Microsoft**

# Azure's data centre infrastructure

Up to **20 data centres** per region
Up to **60 km** of cable between regions

*Region*

Data centres in over **50 regions**
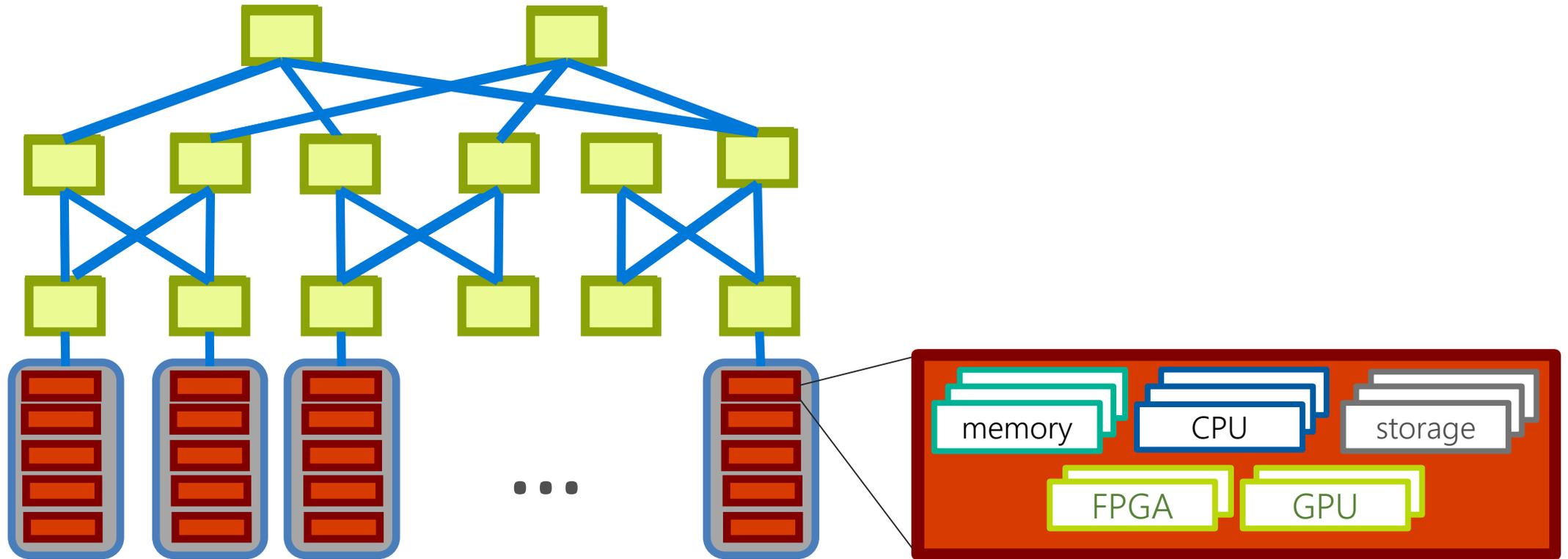~**10,000 km** of fibre-optic cable connecting these

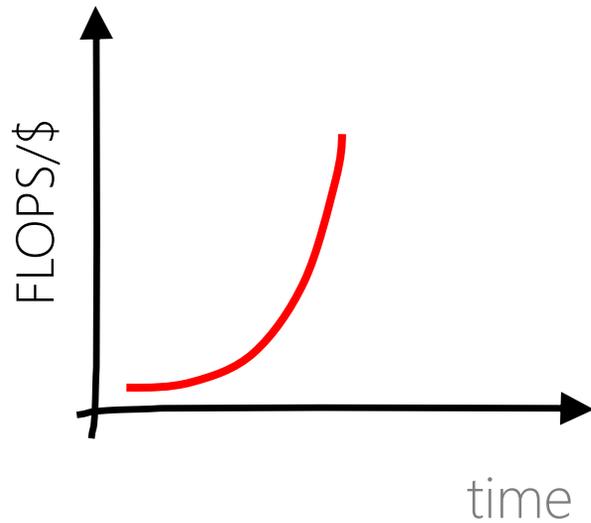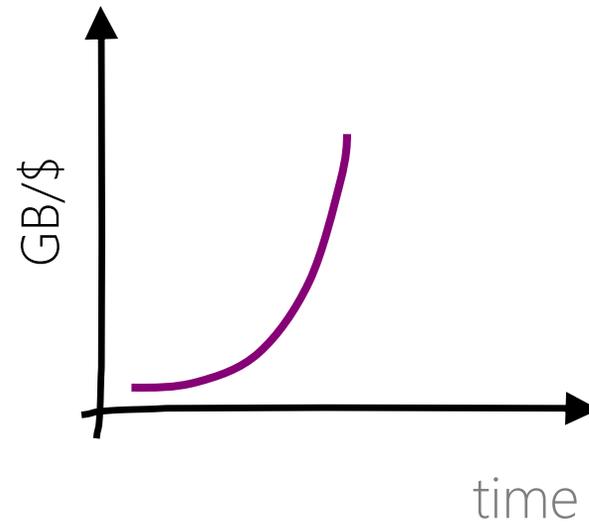~ **100,000 servers** per data centre

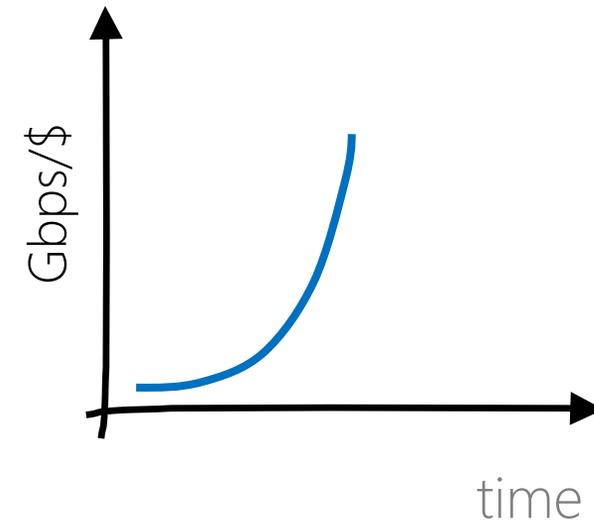# How can we scale the capacity of the cloud?

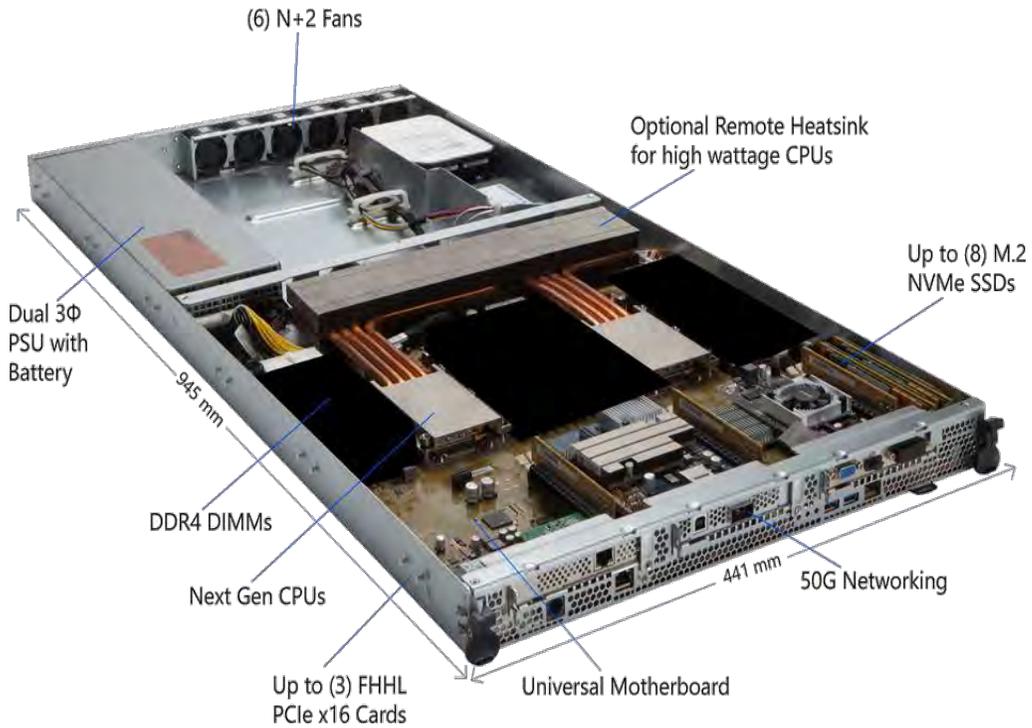We have addressed these challenges in three areas

Compute          Storage          Network

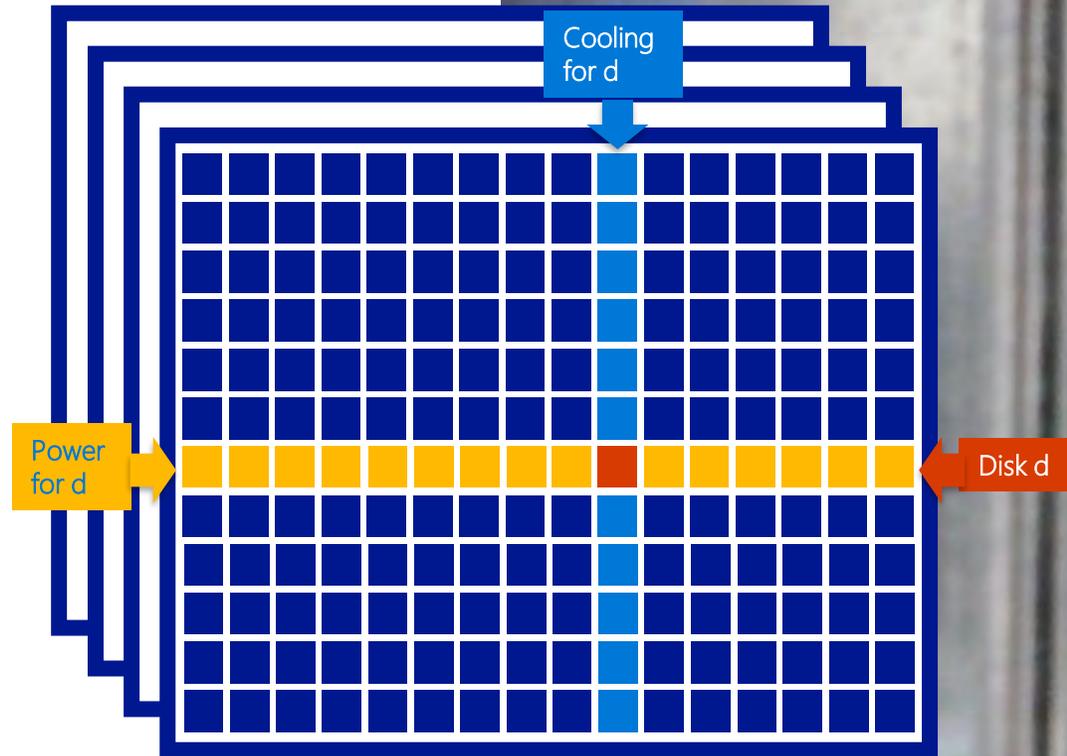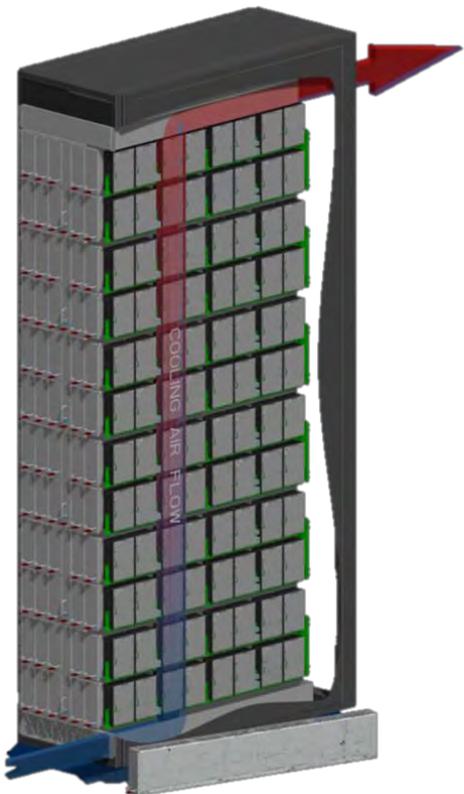# Pushing compute growth

**Project** *Olympus*: hyperscale cloud hardware design in collaboration with OCP (Open Compute Project)

(6) N+2 Fans

Optional Remote Heatsink for high wattage CPUs

Dual 3Φ PSU with Battery

Up to (8) M.2 NVMe SSDs

945 mm

441 mm

DDR4 DIMMs

50G Networking

Next Gen CPUs

Up to (3) FHHL PCIe x16 Cards

Universal Motherboard
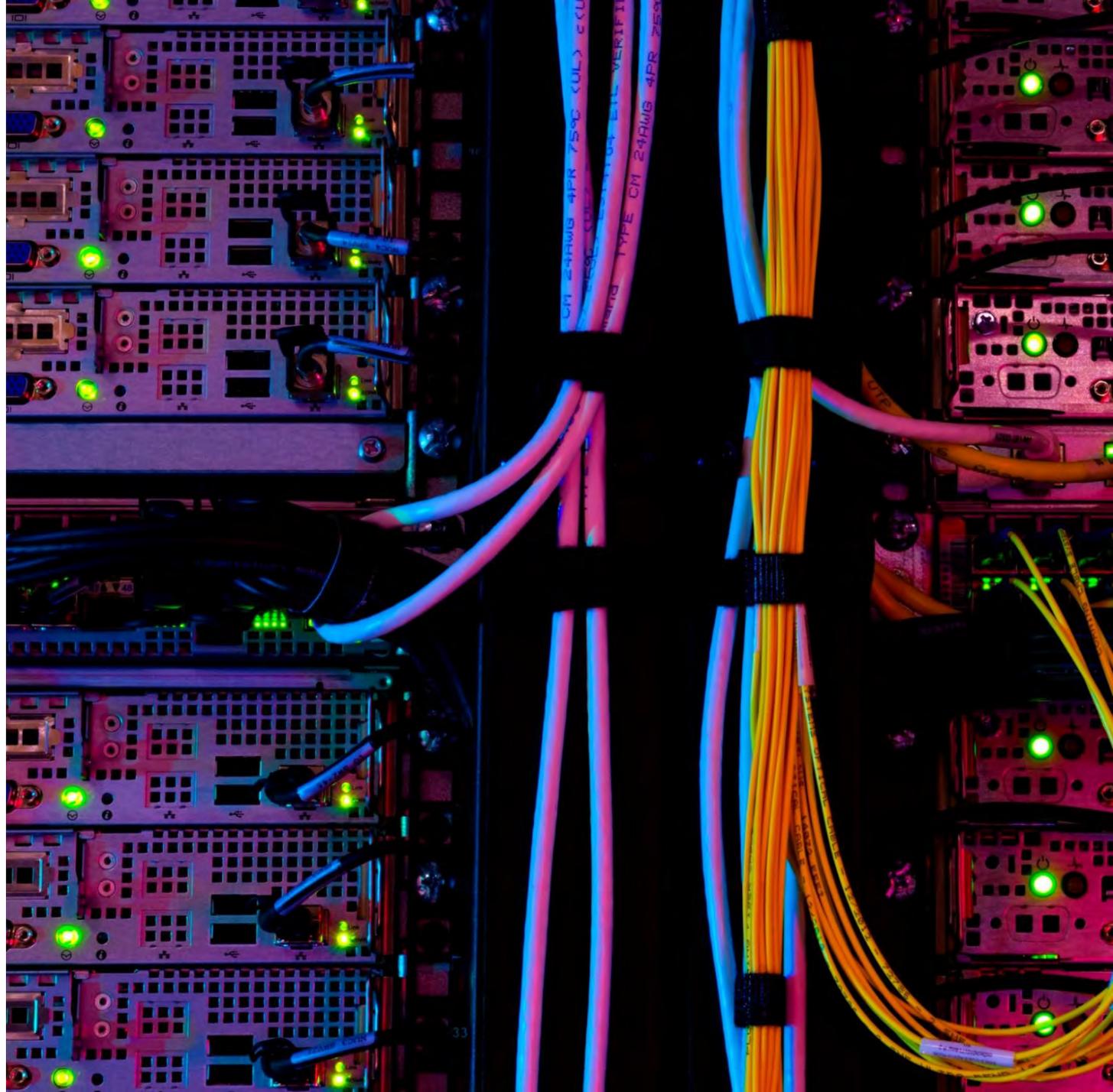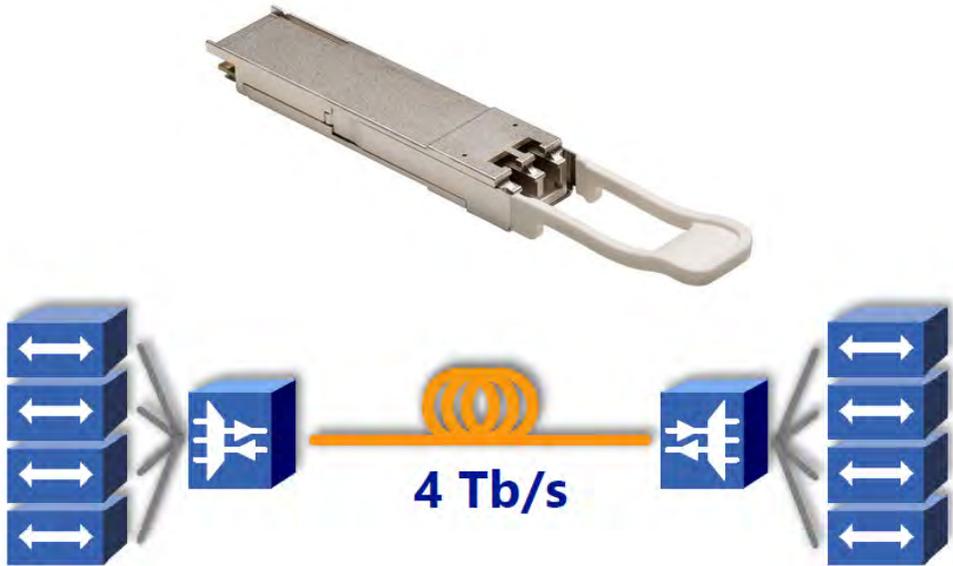
# Pushing storage growth

Project *Pelican*: custom storage racks optimised for cold storage

# Pushing network growth

**Project** *Madison*: custom data centre interconnect technology at multi-terabit speeds

**4 Tb/s**

# We kept up (just)... but the end is nigh... we must respond!

Microsoft

# Impressive advances in optical technologies in the 21ˢᵗ century
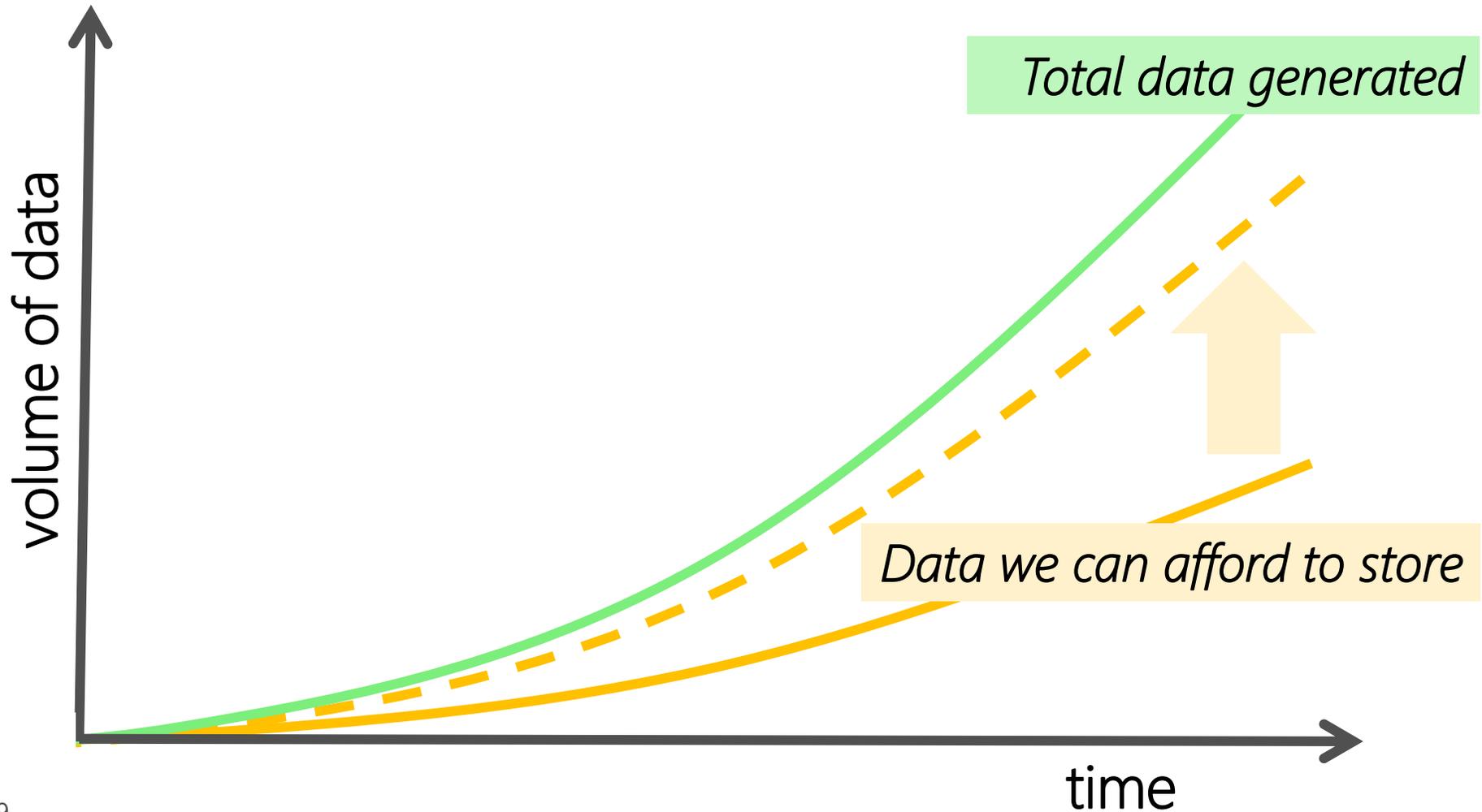
Breakthrough innovation in academic labs across the world:

***Optics for the Cloud*** **Research Alliance**

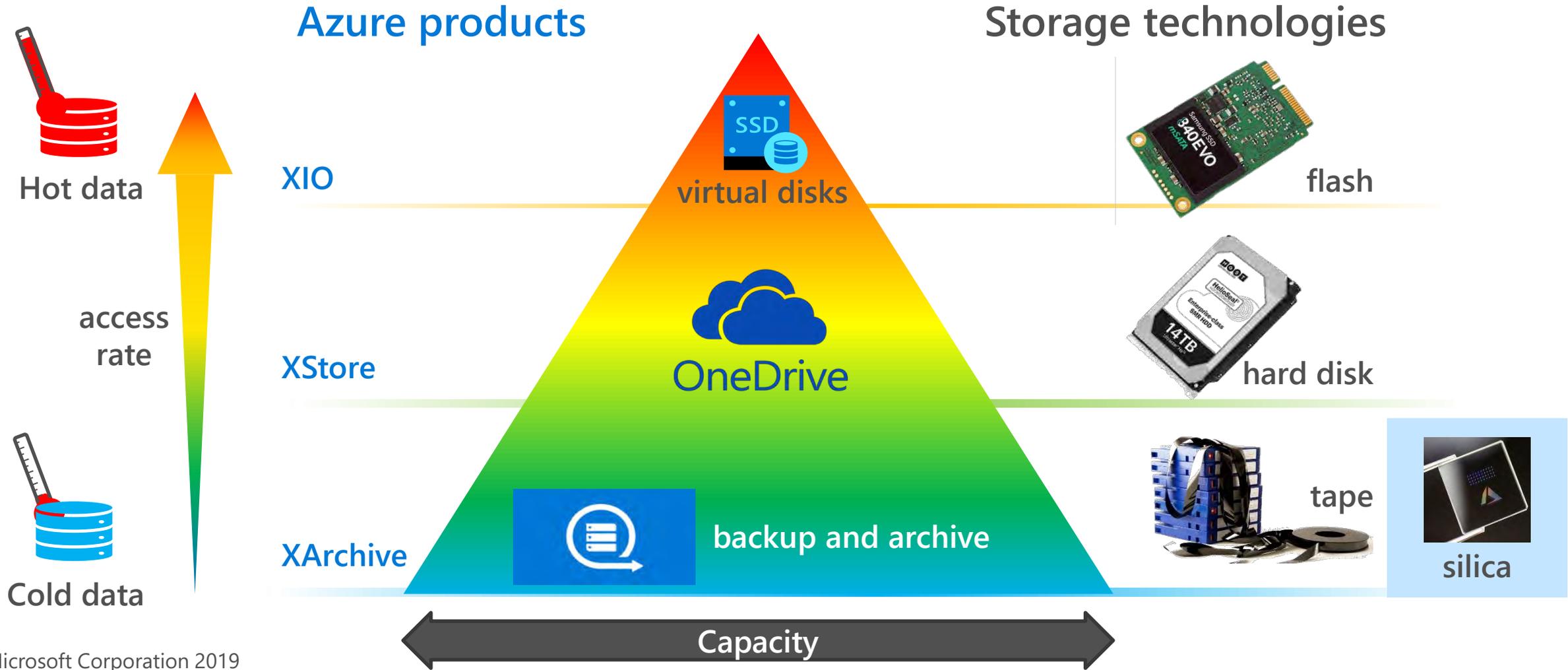**Cambridge, EPFL, Imperial, Southampton, TU/e, UCL, UCSB, etc etc**

Microsoft

Project *Silica*: optical storage for the Cloud

Data volume is growing … we need to reduce $/GB

Total data generated

Data we can afford to store

volume of data

time

© Microsoft Corporation 2019

Today's cloud storage landscape has three tiers

**Azure products**

**Storage technologies**

Hot data

access rate

Cold data

XIO — virtual disks — flash

XStore — OneDrive — hard disk

XArchive — backup and archive — tape, silica

Capacity

© Microsoft Corporation 2019

# Existing archival storage media is magnetic or optical

### hard disk

### magnetic tape

### optical disk

←10 nm

~300 nm

Multiple layers...

# Magnetic media doesn't last

Magnetic media degrades over time

- – Latent sector errors
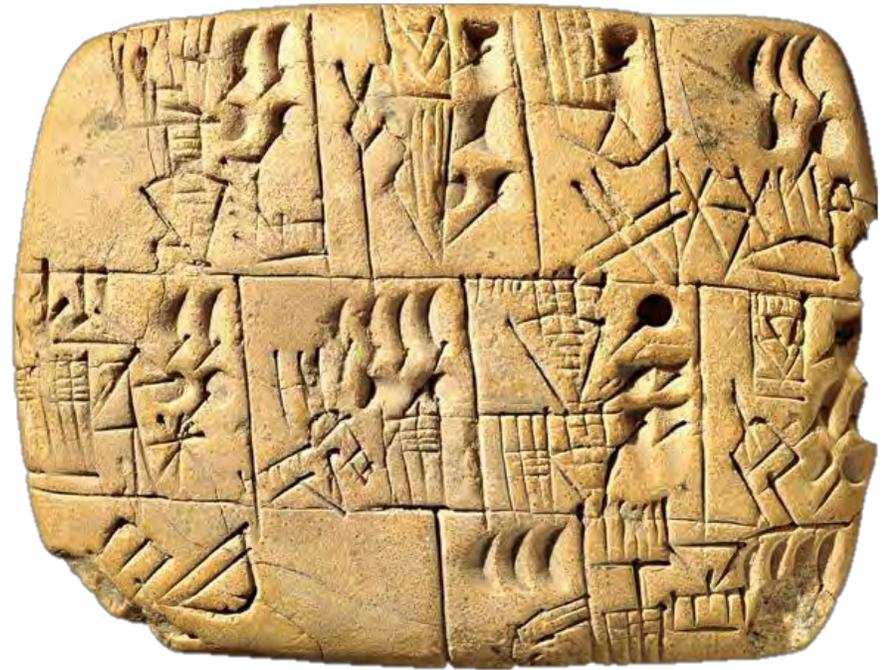
- – Requires scrubbing

Magnetic media has limited life

- – HDD:  3 - 5 years

- – tape:  5 - 7 years

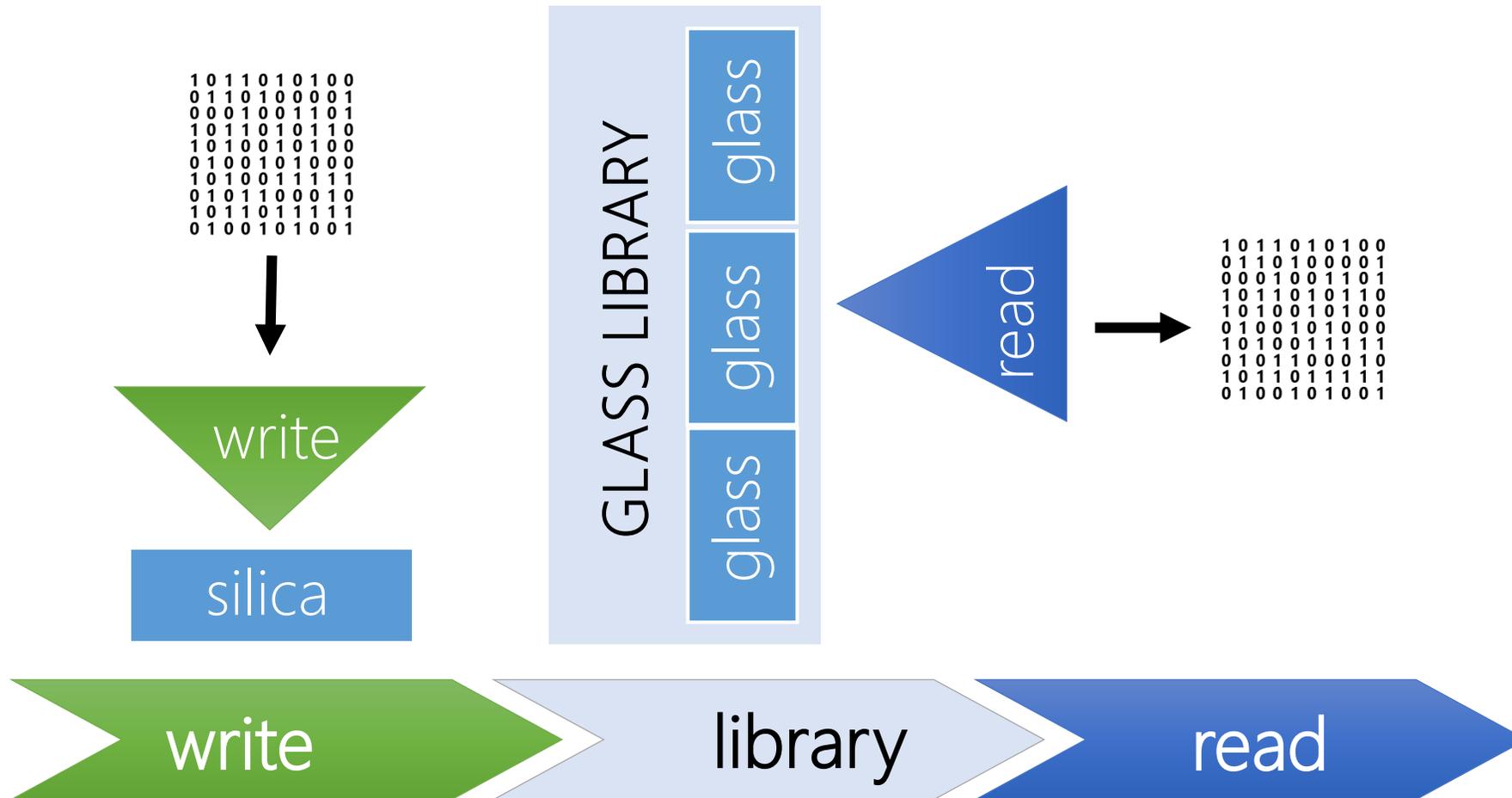# Data can be stored by changing the structure of silica

- Write once, read many
- Persistent – lifetime of millennia
- Immune to EMI/EMP
- No bit rot or disc rot or media decay…
- …hence no need to re-write periodically
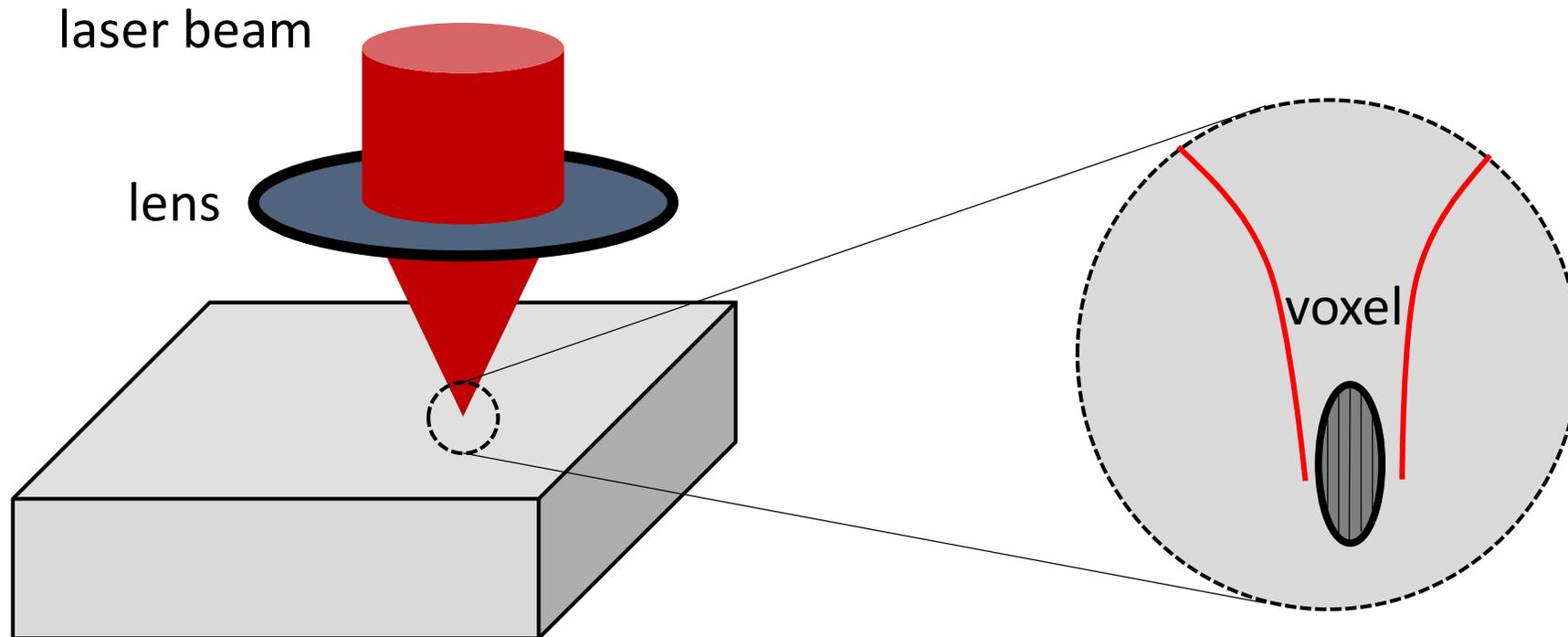- Cheap
- Better access time than tape
- Leave data in place

Cuneiform tablet recording the allocation of beer, 3100-3000 BC.
© Trustees of the British Museum.

Data is written into silica by femtosecond laser pulses

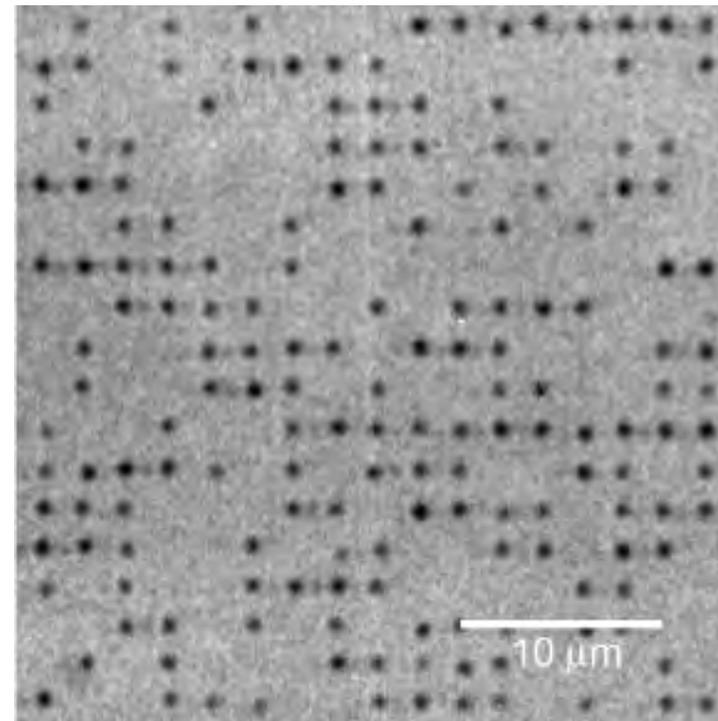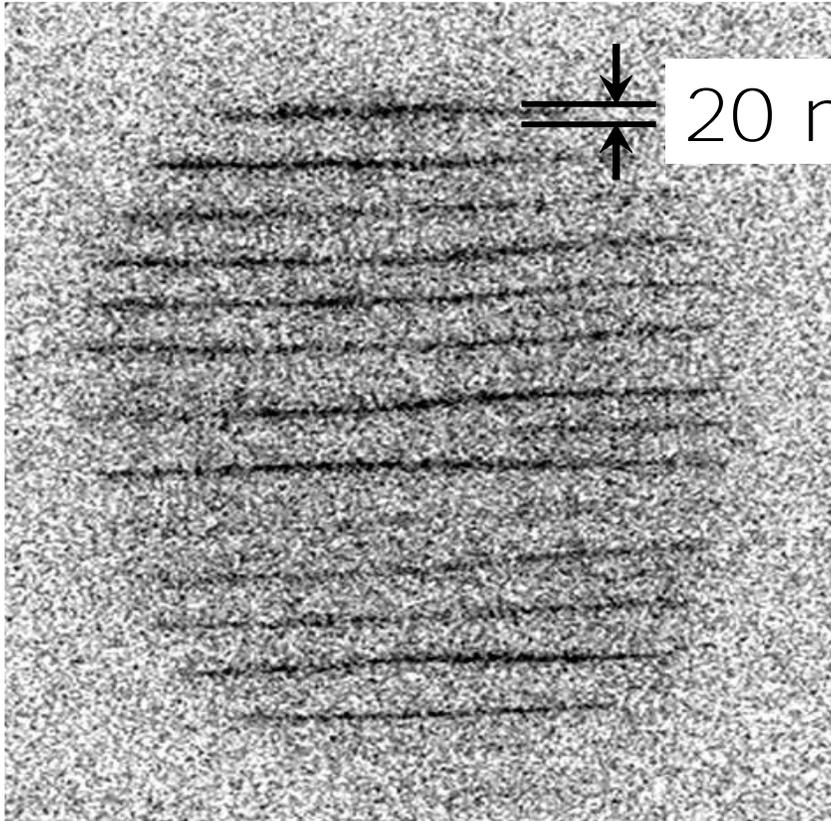laser beam

lens

voxel

# Pulse duration is critical

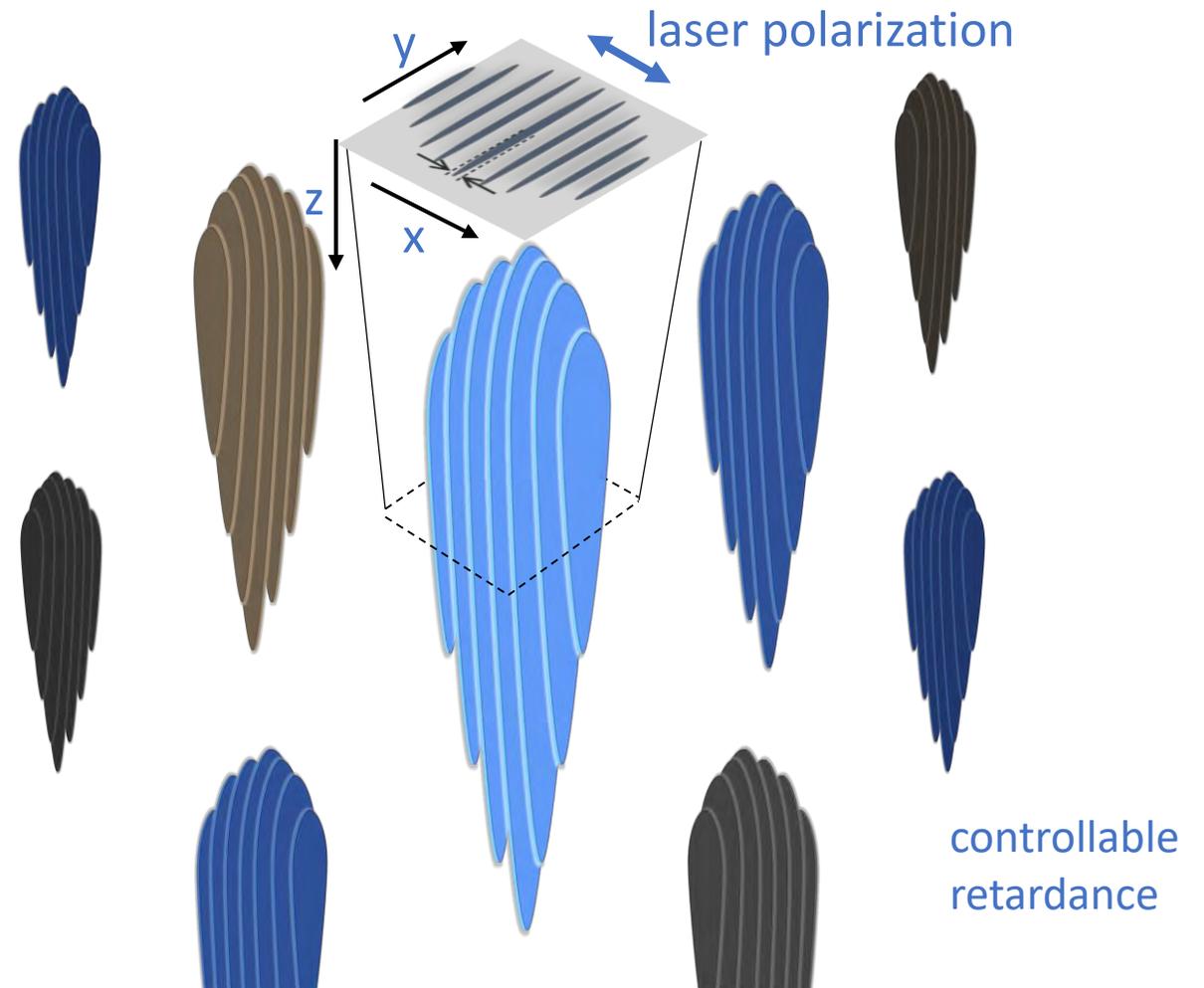picosecond laser induces
*voids with external stress*



femtosecond laser induces
*sub-wavelength nanogratings*

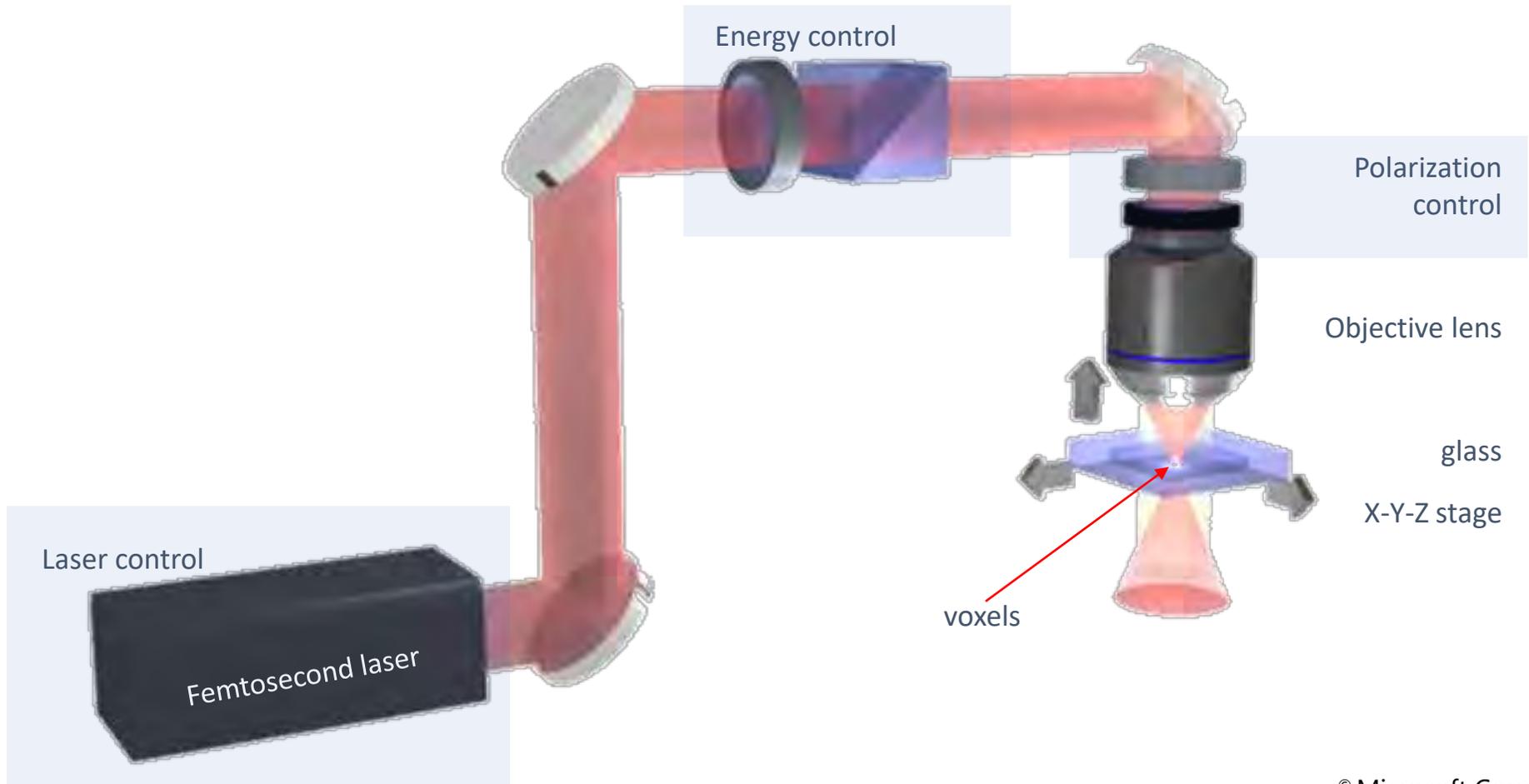# The voxels are sub-wavelength nanogratings

20 nm

*Self-Organized Nanogratings in Glass Irradiated by Ultrashort Light Pulses,* Yasuhiko Shimotsuma, Peter G. Kazansky, Jiarong Qiu, and Kazuoki Hirao Phys. Rev. Lett. **91** (2003)
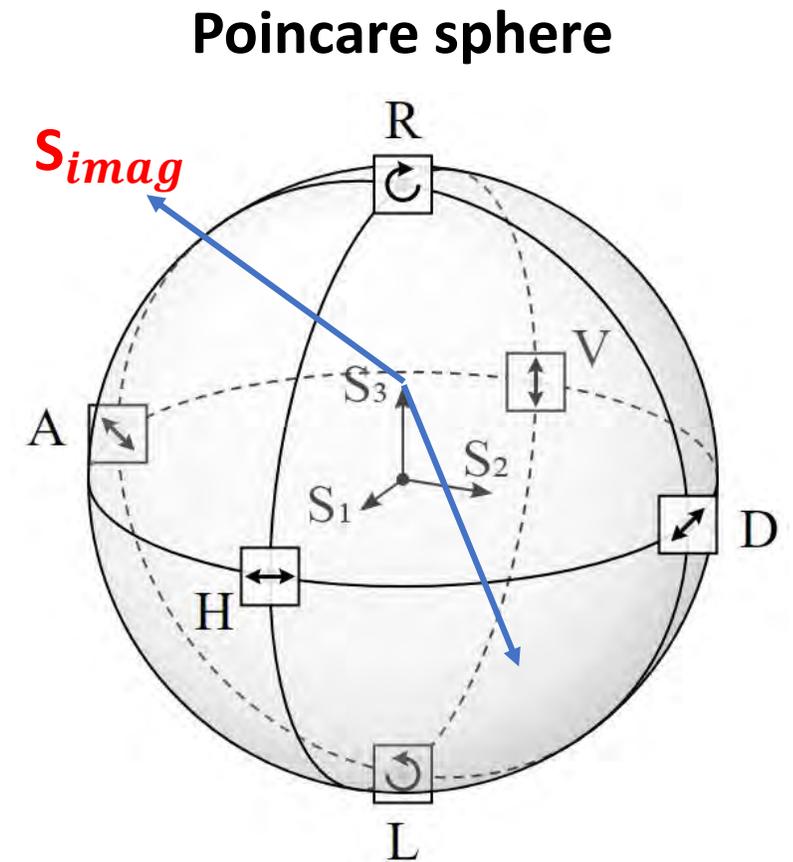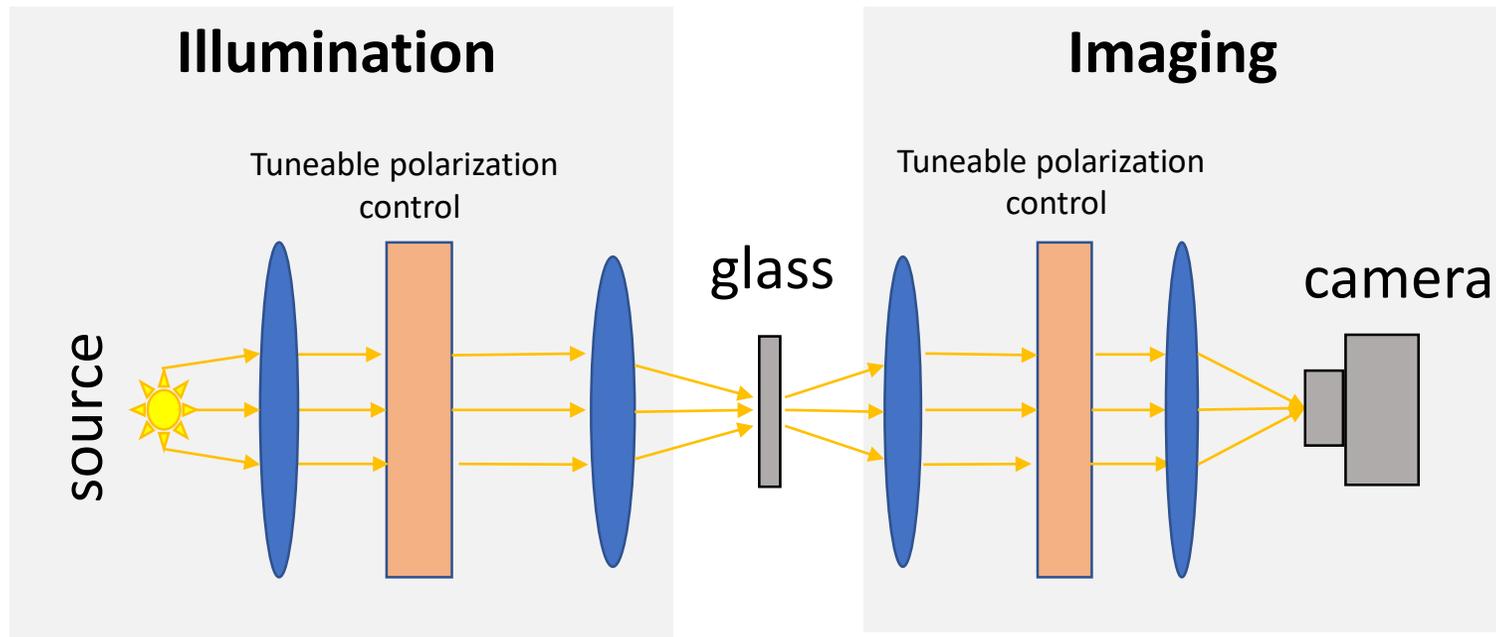
laser polarization

y

z

x

controllable
retardance

# Schematic of WRITE system



Energy control

Polarization control

Objective lens

glass

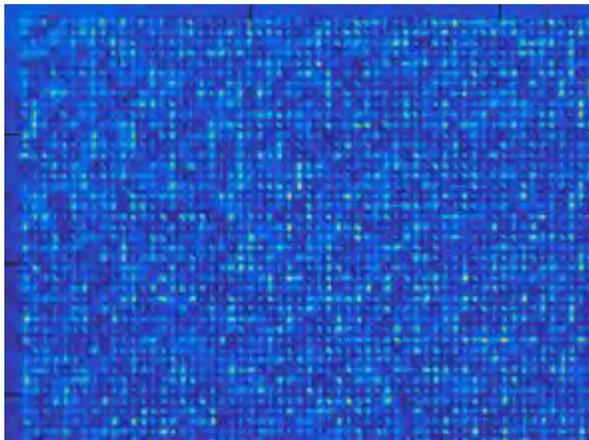X-Y-Z stage

voxels

Laser control

Femtosecond laser

LASER RADIATION

# The READ head is a polarisation microscope

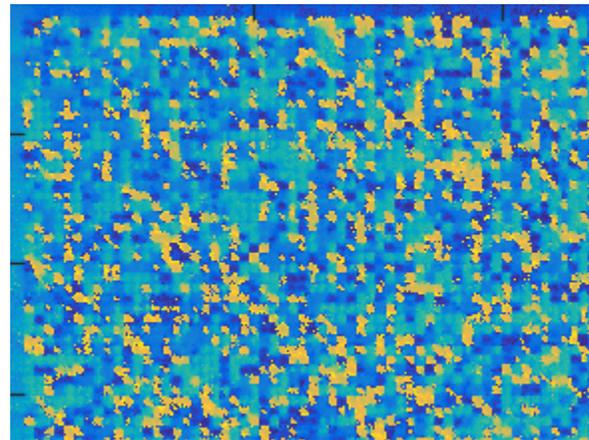- Measure nanogratings under several polarisations of light
- Infer retardance and angle



**Poincare sphere**

# Thank you!

- Tom Empson
- v-toemps@microsoft.com
- aka.ms/silica