



LOCKSS

The Cloudy Outlook for Digital Preservation

Nicholas Taylor ([@nullhandle](#))
Program Manager Emeritus, [LOCKSS](#) and [Web Archiving
Stanford Libraries](#)

[International Conference on Digital Preservation](#)
[Cloud Atlas: Navigating the Cloud for Digital Preservation](#)
17 September 2019

“The Cloud” is playing a growing role in digital preservation.

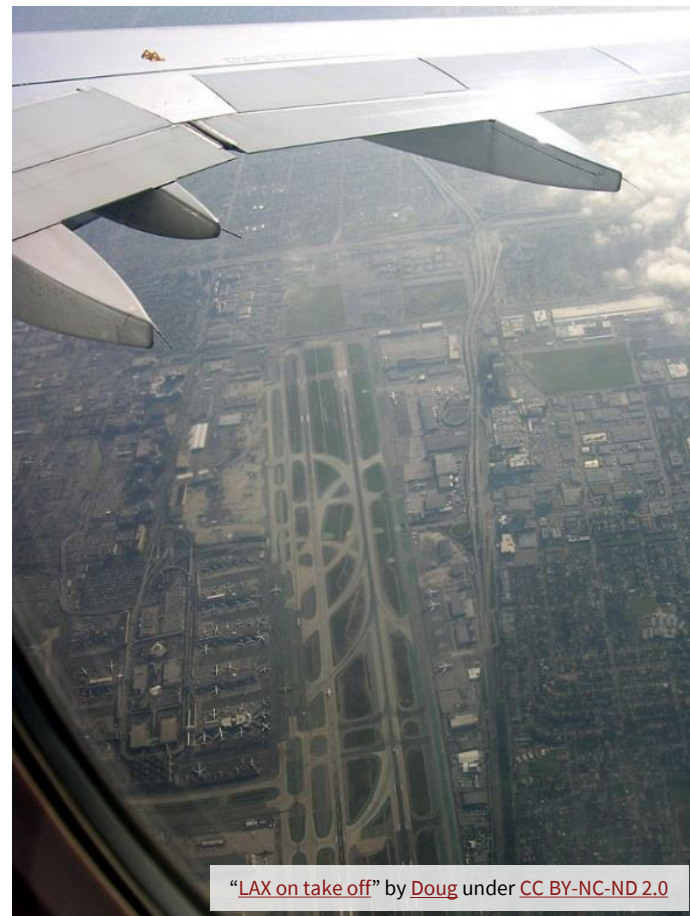
A colorful, stylized illustration of cloud computing concepts. It features a large grey cloud with a red speech bubble saying "OMG!", a yellow speech bubble, and a blue triangle. A red and yellow striped hot air balloon is in the top right. The background is filled with various icons like a smartphone, a laptop, a gear, a lightbulb, a magnifying glass, a document, a network diagram, and a hot air balloon. The bottom of the illustration is a blue wave-like shape containing icons of a laptop, a hot air balloon, a network diagram, and a gear.

Which “The Cloud” we
use, and how we use it,
matters both for our
missions and the likely
success of our efforts.



overview

- threat modeling
- commercial cloud
- community cloud
- wrap up



"LAX on take off" by Doug under CC BY-NC-ND 2.0

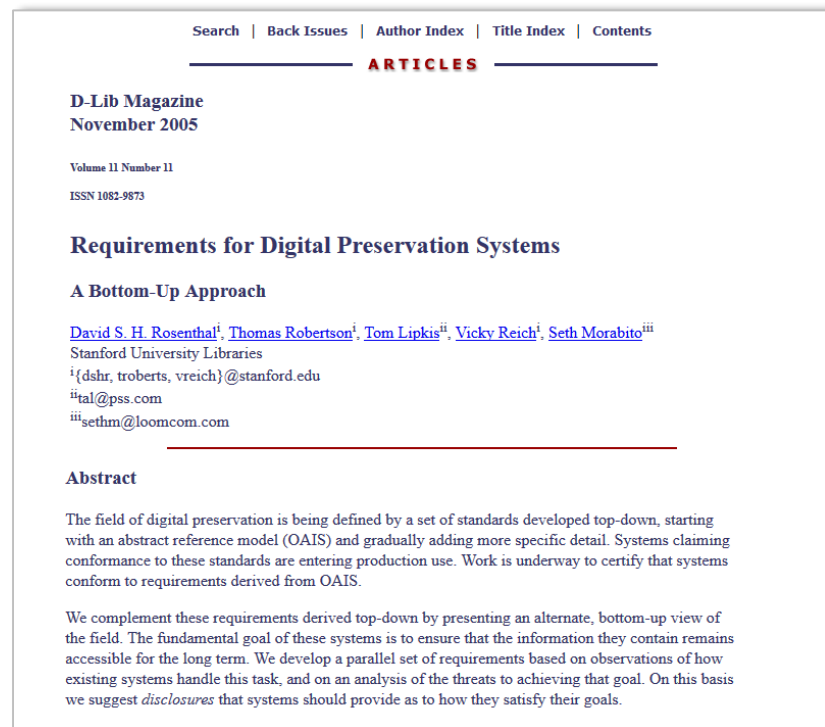


Threat Modeling



threats to digital information

- **failures:** media, hardware, software, network services, economic, organizational
- **obsolescence:** media, hardware, software
- **errors:** communication, operator
- **attacks:** external, internal
- **natural disaster**



[David S.H. Rosenthal et al: "Requirements for Digital Preservation Systems: A Bottom-Up Approach"](#)



community best practices

Table 1: Version 1 of the Levels of Digital Preservation

| | Level 1 (Protect your data) | Level 2 (Know your data) | Level 3 (Monitor your data) | Level 4 (Repair your data) |
|---------------------------------|---|---|---|--|
| Storage and Geographic Location | <ul style="list-style-type: none"> - Two complete copies that are not collocated - For data on heterogeneous media (optical discs, hard drives, etc.) get the content off the medium and into your storage system | <ul style="list-style-type: none"> - At least three complete copies - At least one copy in a different geographic location - Document your storage system(s) and storage media and what you need to use them | <ul style="list-style-type: none"> - At least one copy in a geographic location with a different disaster threat - Obsolescence monitoring process for your storage system(s) and media | <ul style="list-style-type: none"> - At least three copies in geographic locations with different disaster threats - Have a comprehensive plan in place that will keep files and metadata on currently accessible media or systems |
| File Fixity and Data Integrity | <ul style="list-style-type: none"> - Check file fixity on ingest if it has been provided with the content - Create fixity info if it wasn't provided with the content | <ul style="list-style-type: none"> - Check fixity on all ingests - Use write-blockers when working with original media - Virus-check high risk content | <ul style="list-style-type: none"> - Check fixity of content at fixed intervals - Maintain logs of fixity info; supply audit on demand - Ability to detect corrupt data - Virus-check all content | <ul style="list-style-type: none"> - Check fixity of all content in response to specific events or activities - Ability to replace/repair corrupted data - Ensure no one person has write access to all copies |

NDSA Levels of Preservation Working Group : “NDSA Levels of Digital Preservation”



what threats are addressed?

- **failures:** media, hardware, software, network services, economic, organizational
- **obsolescence:** media, hardware, software
- **errors:** communication, operator
- **attacks:** external, internal
- **natural disaster**

Table 1: Version 1 of the Levels of Digital Preservation

| | Level 4 (Repair your data) |
|---------------------------------|---|
| Storage and Geographic Location | <ul style="list-style-type: none">- At least three copies in geographic locations with different disaster threats- Have a comprehensive plan in place that will keep files and metadata on currently accessible media or systems |
| File Fixity and Data Integrity | <ul style="list-style-type: none">- Check fixity of all content in response to specific events or activities- Ability to replace/repair corrupted data- Ensure no one person has write access to all copies |

NDSA Levels of Preservation Working Group : “[NDSA Levels of Digital Preservation](#)”



what threats are discounted?

- **failures:** media, hardware, software, network services, economic, organizational
- **obsolescence:** media, hardware, software
- **errors:** communication, operator
- **attacks:** external, internal
- **natural disaster**

Table 1: Version 1 of the Levels of Digital Preservation

| | Level 4 (Repair your data) |
|---------------------------------|---|
| Storage and Geographic Location | <ul style="list-style-type: none">- At least three copies in geographic locations with different disaster threats- Have a comprehensive plan in place that will keep files and metadata on currently accessible media or systems |
| File Fixity and Data Integrity | <ul style="list-style-type: none">- Check fixity of all content in response to specific events or activities- Ability to replace/repair corrupted data- Ensure no one person has write access to all copies |

NDSA Levels of Preservation Working Group : “NDSA Levels of Digital Preservation”



Amazon cloud crash wipes out customer data; Will users be compensated?

By **Molly McHugh** —
Posted on April 28, 2011 10:57 am



SECURITY

Ohio's Recent Spate of Cyberattacks Is Indicative of the National Trend



FOX NEWS

Home Video Politics U.S. Opinion Business

Angry Employee Deletes All of Company's Data

Published January 24, 2008

Fox News



Call it a tale of revenge gone wrong.



Yahoo Quietly Pulls The Plug On Geocities

Posted Apr 23, 2009 by **Leena Rao (@leenarao)**

Sorry, new GeoCities accounts are no longer available.

Current GeoCities customers:

After careful consideration, we have decided to close GeoCities later this year. We'll have more details this summer. For now, please sign in or visit the help center for more information.

• Sign in to GeoCities

Save 50% on a web site with Yahoo! Web Hosting.

- Get a personalized web address (e.g., www.lean2yoga.com).
- Design a great-looking site with easy-to-use tools.
- Never worry about ads appearing on your site.
- Get unlimited disk space to store all your pictures and audio tracks.
- Welcome as many visitors as you like, thanks to unlimited data transfer. Unlimited.
- Get answers even at 2 a.m., with 24-hour phone support.

[Learn more](#) | [View demo](#)



FUTURE PERFECT EXPLAINERS MORE

"Climate change" and "global warming" are disappearing from government websites

The deletions follow a pattern of policy changes on climate change under the Trump administration.

By Umair Irfan | Updated Jan 11, 2018, 12:30pm EST



REUTERS

Davos

The Trump Effect

Politics

#SCIENCE NEWS

JULY 20, 2009 / 3:19 PM / 9 YEARS AGO

Moon landing tapes got erased, NASA admits

Maggie Fox, Health, Science Editor



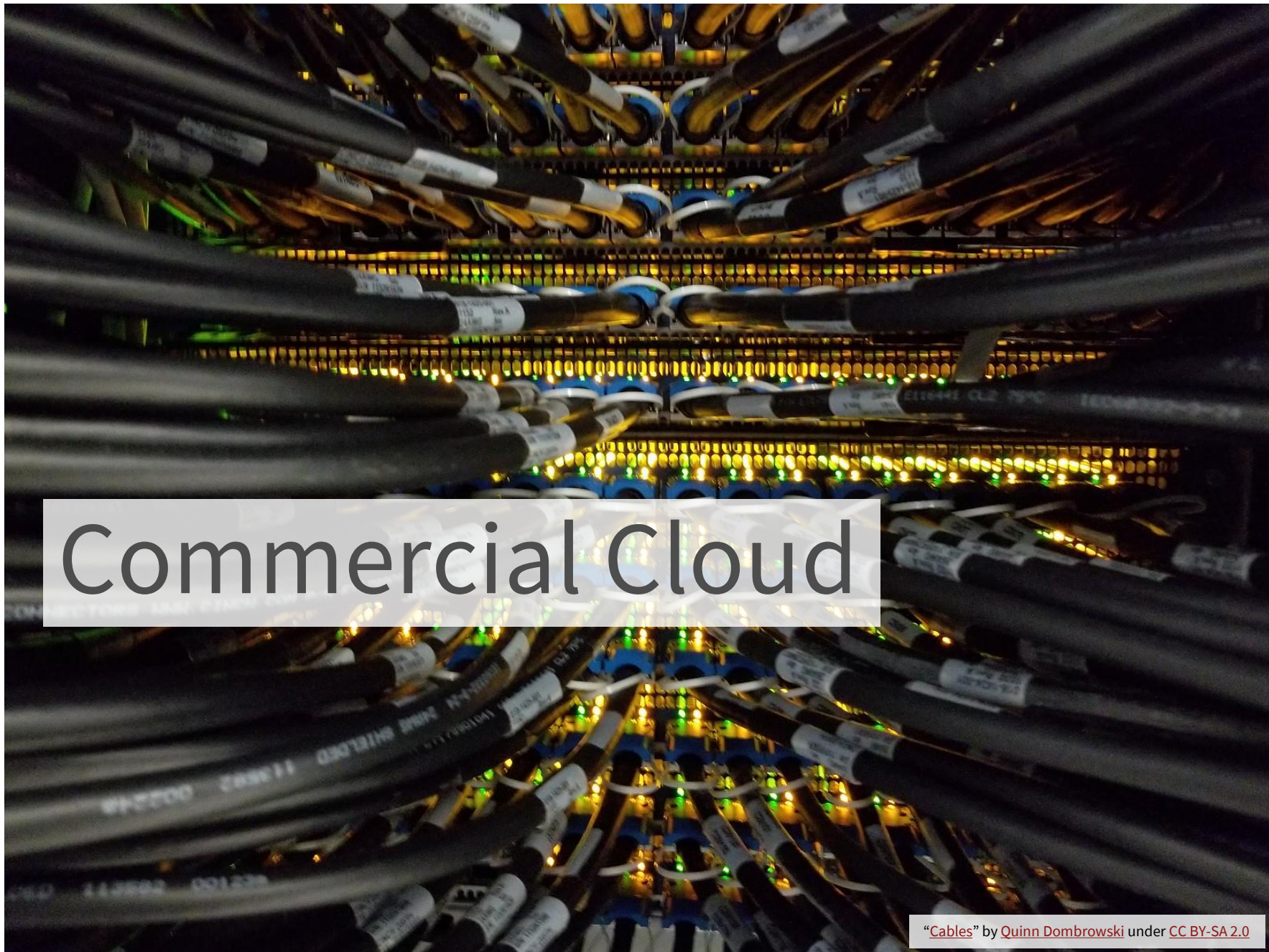
WASHINGTON (Reuters) - The original recordings of the first humans landing on the moon 40 years ago were erased and re-used,



understand + mitigate threats

- long-term **data integrity** is hard
- needs architecture informed by **actual leading threats** to data
- don't underestimate:
 - people making **mistakes**
 - **attacks** on information
 - organizational **failure**







commercial cloud considerations

- on-demand access
- economic lock-in
- reliability caveats
- opaque data integrity
- security configuration
- non-operational externalities





on-demand access

- minimize **idle built-out infrastructure** needed for:
 - long-tail access
 - data integrity checks
- shift costs for **corpus-level use cases**
- but **metered access**:
 - complicates cost modeling
 - works less well for popular or oft-accessed content





economic lock-in

- CapEx → OpEx yields **financial flexibility**
- financial flexibility **less valuable for inflexible commitment** (i.e., long-term preservation)
- prices decline but maybe **not as quickly** as in competitive market for local hardware refresh
- mitigating strategy: maintain **a local copy**

| long-term storage services (1 PB, 1 month) | | | | |
|--|---------|---------|----------|----------------|
| service | ingest | store | export | lock-in factor |
| AMZN Glacier | \$2,250 | \$4,000 | \$55,240 | 13.8x |
| GOOG Coldline | \$3,600 | \$7,000 | \$83,860 | 12x |
| MSFT Archive | \$6,350 | \$2,000 | \$16,260 | 8.1x |

David S.H. Rosenthal : “[Cloud for Preservation](#)”



“Dam II” by [Craig Bennett](#) under [CC BY-NC 2.0](#)



reliability caveats

- “**11 nines**” of reliability?
 - modeled on hardware failure
 - accounts for $\frac{1}{3}$ of data losses
- $\frac{2}{3}$ of data losses due to **less rationalizable factors**: attacks, errors, software failures
 - highly centralized infrastructure more vulnerable
- **chance of billing error** interrupting service non-trivially more significant than risk of loss suggested by reliability estimate



“Frayed, but holding up” by [Edna Winti](#) under [CC BY 2.0](#)



opaque data integrity

- **feature**, not a bug?
- hashing data in situ **requires trusting** that the service has performed computation rather than reporting cached value
- may be **prohibitively expensive** to retrieve content to a trusted environment to perform hashing





security configuration

- **monoculture** vulnerabilities
- greater affordances, better defaults for **on-premise security**
- consistent leaks from **misconfigured cloud services** suggest security is a challenge





The Telegraph Log in

News Politics Sport Business Mon

See all Tech

Technology Intelligence

Tech giants face tax avoidance crackdown

share Save 3

GIZMODO

LATEST REVIEWS SCIENCE 109 FIELD G



GOODBYE BIG FIVE

I Cut the 'Big Five' Tech Giants From My Life. It Was Hell

Share Tweet

Home Mail News Finance Spo

YAHOO! FINANCE Search for news

Finance Home Watchlists My Portfolio

Exclusive: 'Ready to stomp on it': Documents reveal staggering power of tech giant lobbying

CPO MAGAZINE

DATA PRIVACY INSIGHTS

5 MIN READ

Is Big Tech Too Big to Fail?

NICOLE LINDSEY · MAY 28, 2018

THE VERGE

TECH SCIENCE ENTERTAINMENT MORE

PODCASTS THE VERGECAST TECH 21

Big Tech's problem is its lack of competition

YaleEnvironment360

Explore Search About E360




ILLUSTRATION BY MATT ROTA

Energy Hogs: Can World's Huge Data Centers Be Made More Efficient?



Community Cloud

"The Ardent Mobile Cloud Platform rains on the DPW Parade, Burning Man 2013" by Neil Girling under [CC BY-NC-ND 2.0](#)



not all clouds the same

- “*The Cloud is just somebody else’s computer*”
- values-aligned partnerships to **build private clouds** e.g.,
 - consortial/community
 - focused on particular content types (e.g., software, web archives)
 - for computational research





community cloud considerations

- sustain community capacity
- flexibility + interoperability
- diversity + risk mitigation
- pilot models





sustain community capacity

- can we still claim to have **custody + intellectual control** over content stored in commercial cloud?
- can we afford to outsource functions **core to mission** to commercial cloud?
- **scholarly publishing** is an example of a service ceded to commercial providers





flexibility + interoperability

*“[T]he needs of today’s diverse scholarly communities are not being met by the existing largely uncoordinated scholarly infrastructure, which is dominated by vendor products that take **ownership of the scholarly process and data**. We intend to create **a new open infrastructure system** that will enable us to work in a more integrated, collaborative and strategic way. It will support **global connections and consistency** where it is appropriate, and **local and contextual requirements** where that is needed.”*

[Invest in Open Leadership: “Preamble, The Why”](#)



diversity + risk mitigation

- lots of copies is **necessary but not sufficient**
- central points of failure can **undermine all copies at once**
- **multi-organizational** preservation storage provides:
 - **resilience** against organizational failure
 - **diversity** in technical infrastructure



["Domino's"](#) by [david pacey](#) under [CC BY 2.0](#)



pilot models

- in original, Global LOCKSS Network, **all nodes stored copies**
- private LOCKSS networks moving towards **hosted service models**
- **subset of institutions** host infrastructure w/ governance by + funding from broader community
- **Stanford + trusted partners** may serve as anchor storage hosts

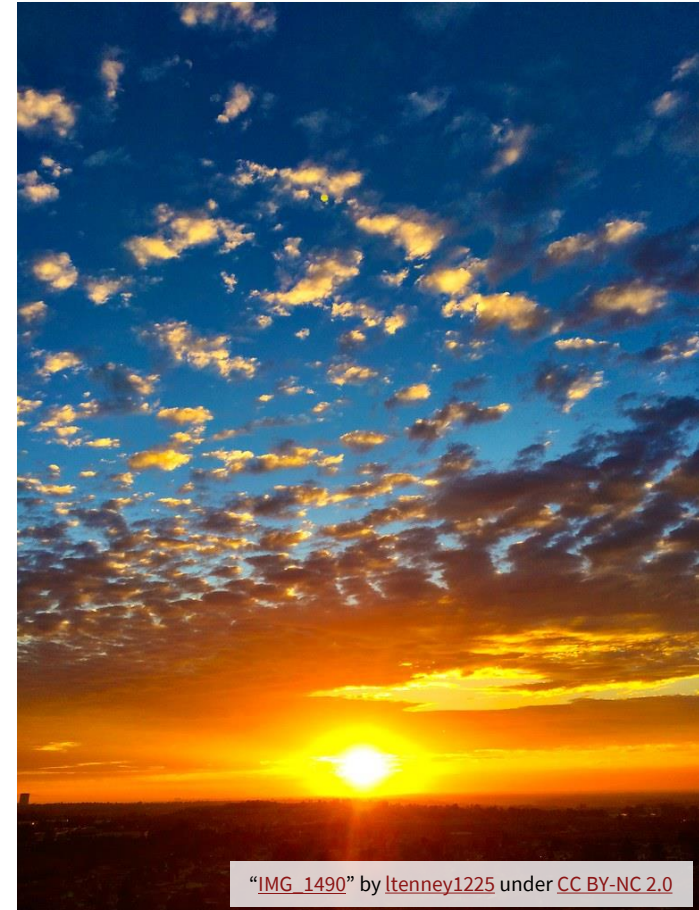




Wrap Up

deliberate cloud strategy

- let's be cautious about **reducing future flexibility**
- let's understand the **meaningful differences** between use cases
- let's be mindful of **trade-offs**
- let's consider what else we can do on **open infrastructure, together**





Questions

[“Any Questions?”](#) by [Matthias Ripp](#) under [CC BY 2.0](#)