



DIGITAL PRESERVATION AT THE NATIONAL LIBRARY OF NORWAY

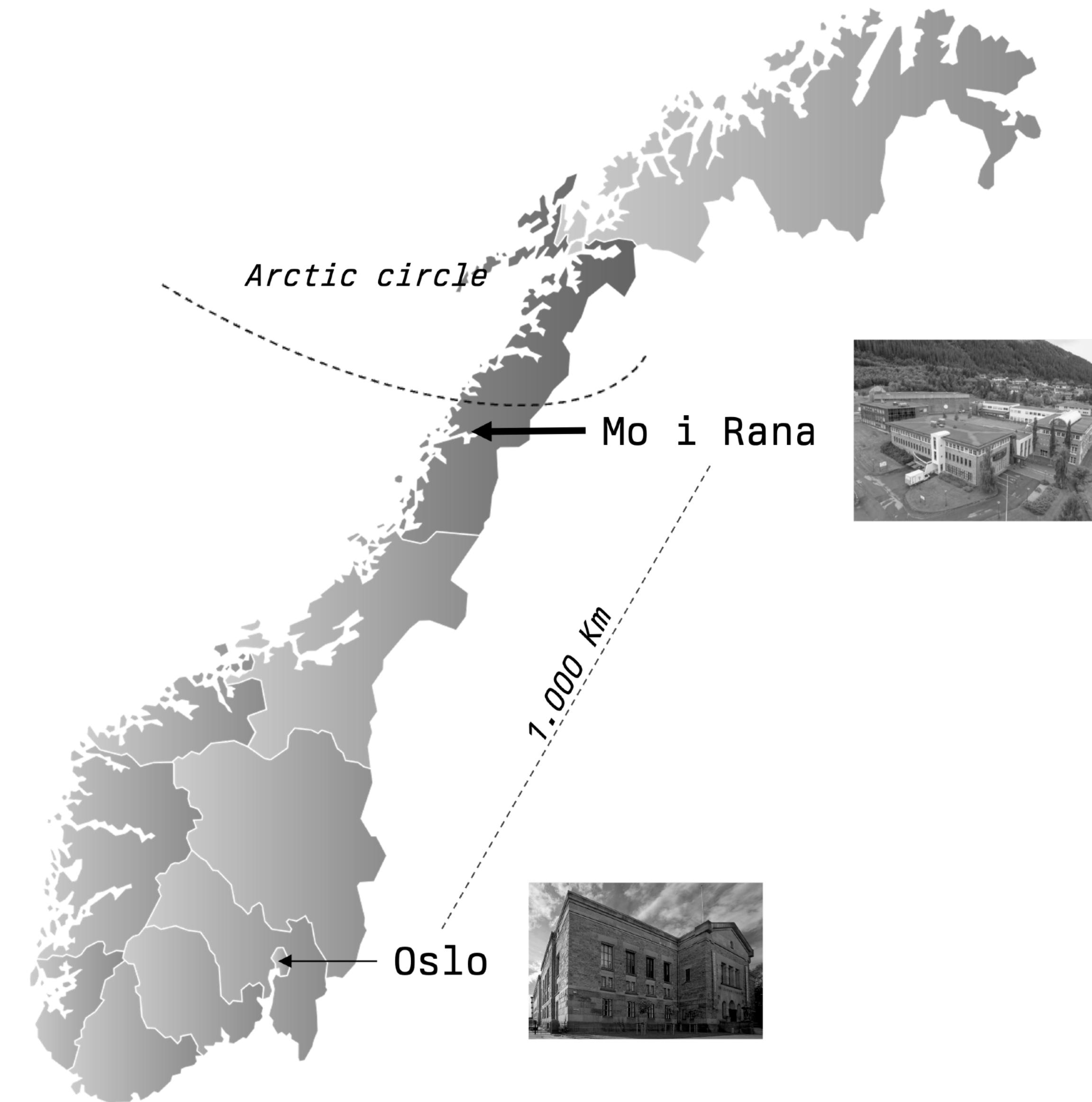
Torbjørn Pedersen 2024-29-05

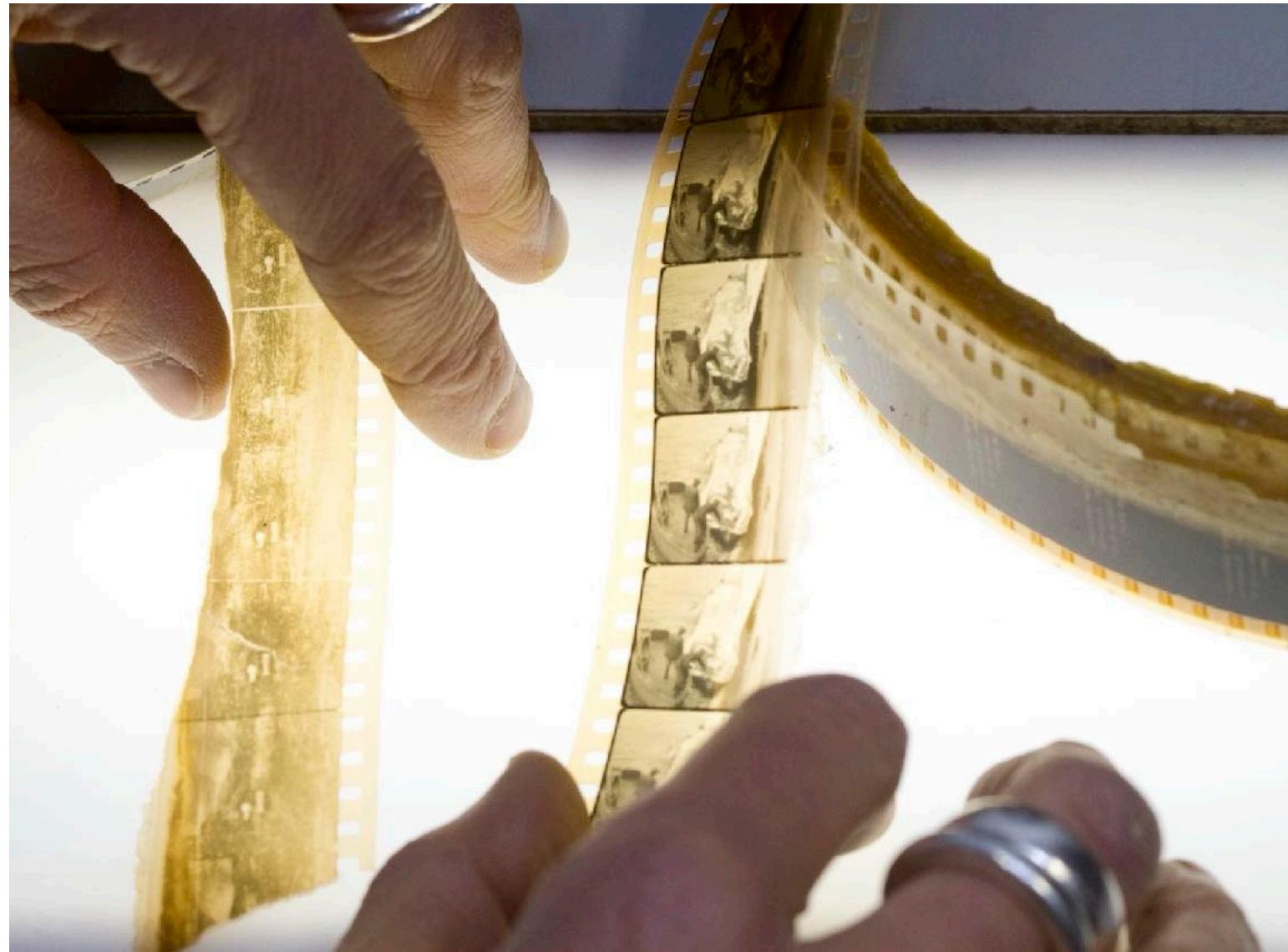
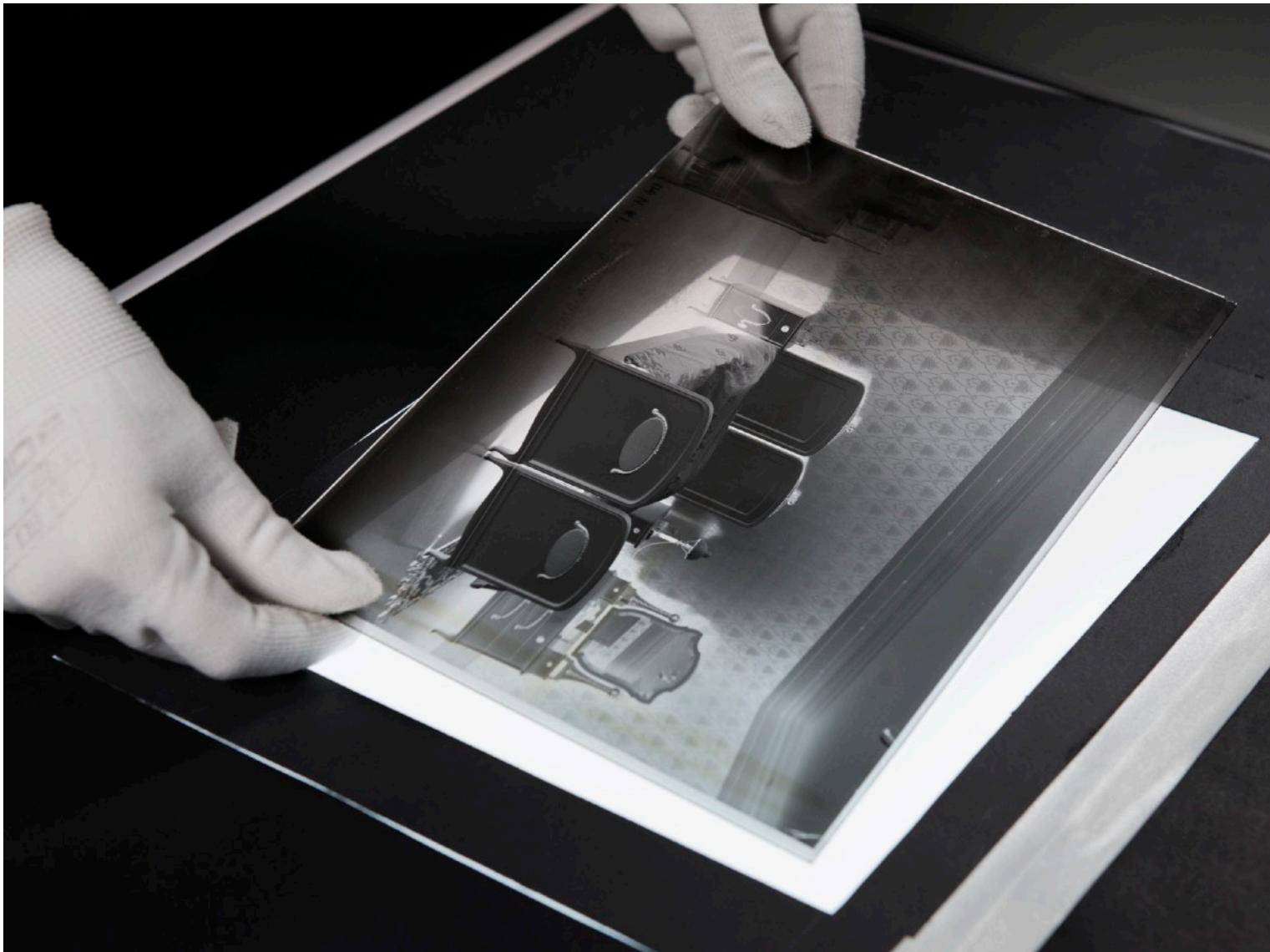
Past → present → future

From digital storage
→
to digital preservation

THE NATIONAL LIBRARY OF NORWAY (NLN)

- ▶ Almost 600 employees at 2 locations: Oslo and Mo i Rana
- ▶ **Governed by:**
 - ▶ The Legal Deposit act
- ▶ **Responsible for:**
 - ▶ Collecting, preserving, making available all content published in Norway
 - ▶ (also historical collections across all media types)





DIGITISATION PRE-2006

Early beginnings/Ancient history

- ▶ Large scale digitisation (at the time)
 - ▶ Photography
 - ▶ Audio and radio
- ▶ Formed basis for what was to come... (2006→)
- ▶ Basic bit storage and backups
 - ▶ Two data loss incidents

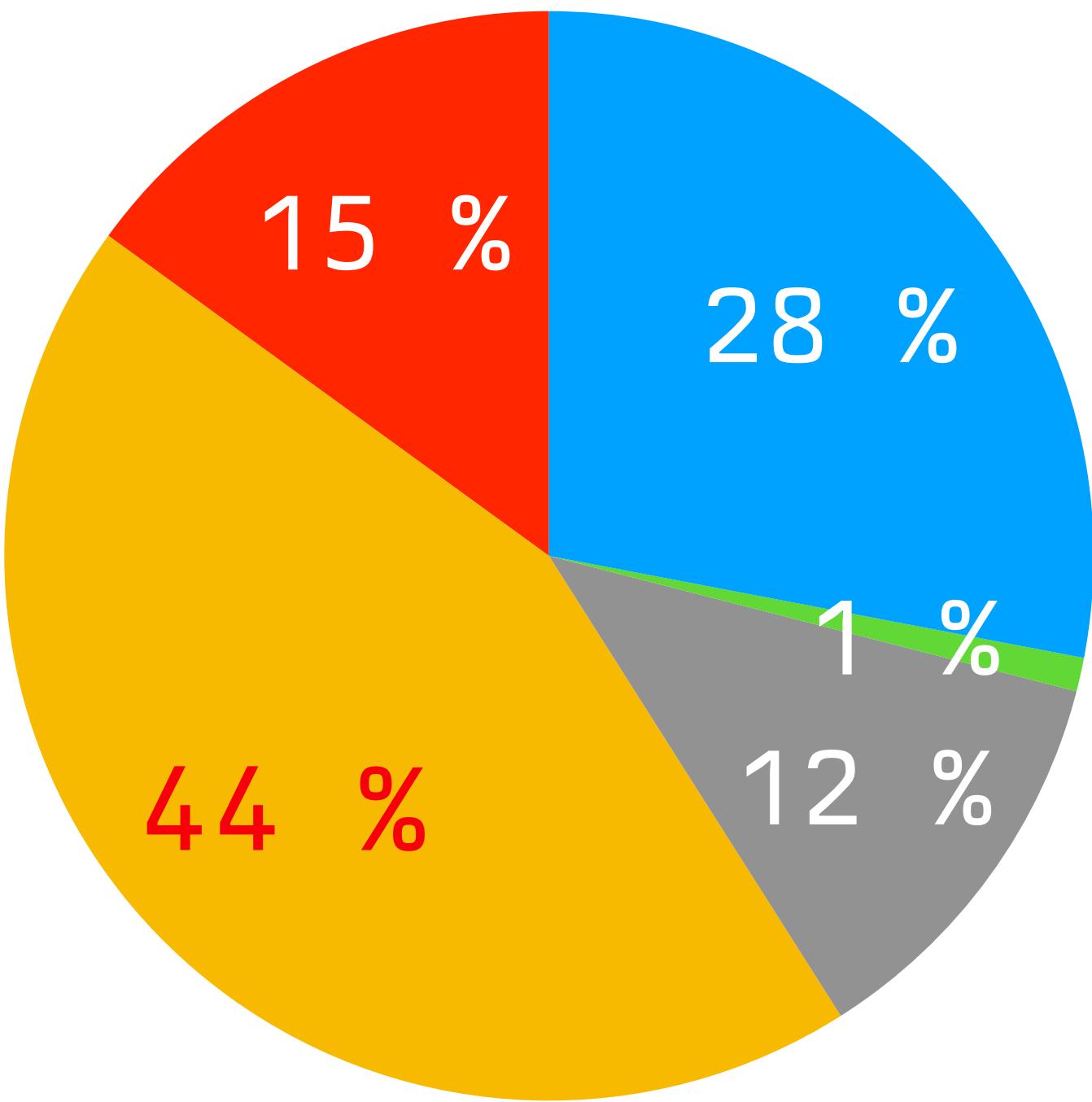
MASS DIGITIZATION 2006→2022

- ▶ **New mantra:**
 - ▶ "If it's not online, it doesn't exist"
 - ▶ "Analog objects wear out when accessed, digital objects don't"
- ▶ **Decision:** digitize the entirety of our analog collection for preservation (and access)
- ▶ **Mass digitisation** efforts for all media types
 - ▶ "Production lines" established - investment in people, equipment, automation
 - ▶ 30+ different production lines created over time
- ▶ Born digital deposits gradually increasing

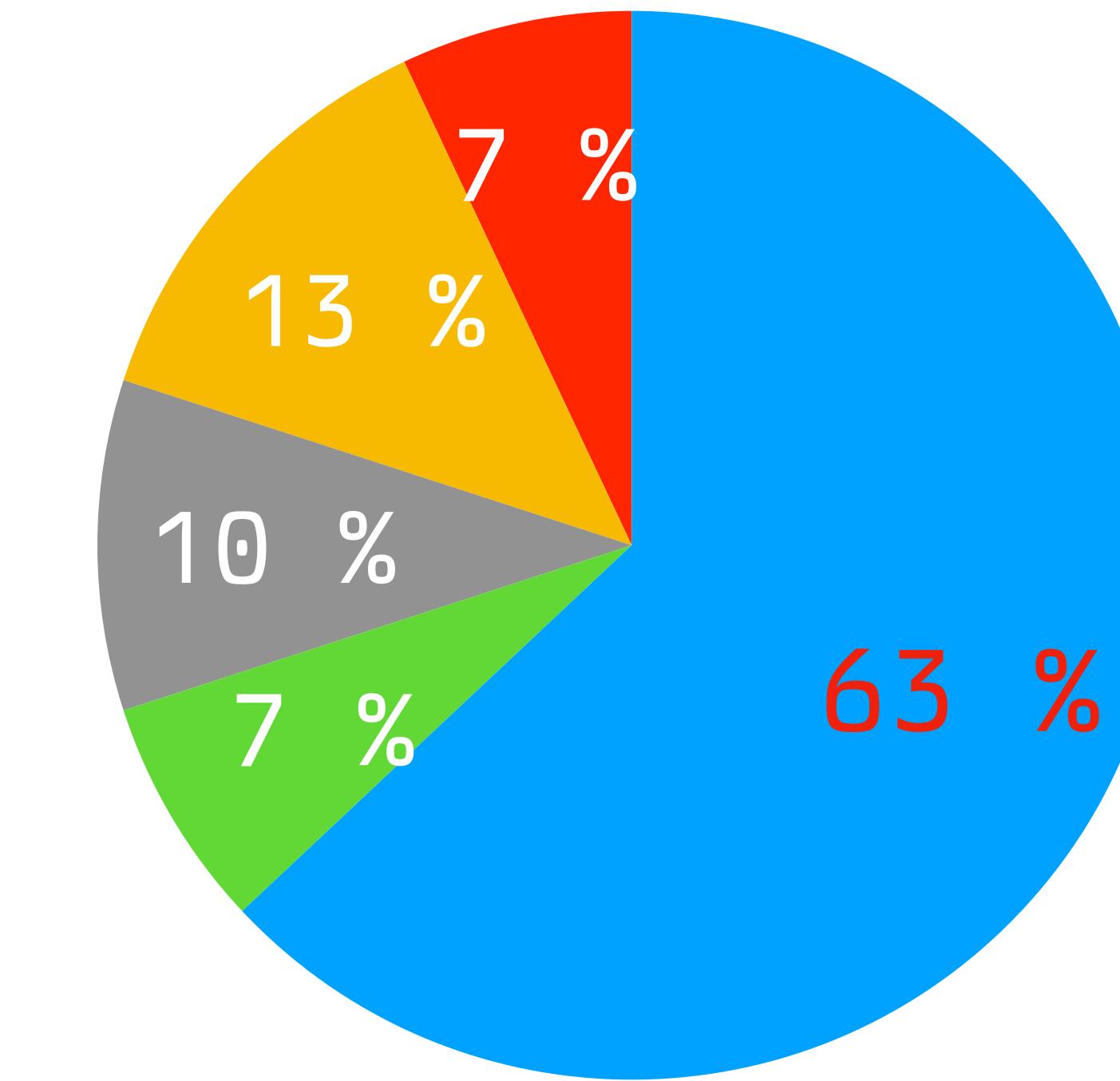
DIGITAL PRESERVATION 2006→2022

- ▶ **3-2-1 storage policy (bit preservation)**
 - ▶ 3 copies (disk+tape+tape)
 - ▶ 2 technologies: disk (luxurious!) and tape
 - ▶ 1 copy at a “off-site” location
- ▶ **New bit repository (2007): SAM-FS (Oracle HSM)**
 - ▶ Solaris OS
 - ▶ File based storage
 - ▶ Mix of different disk and tape technologies (no vendor lockin)
 - ▶ **Disk:** SUN/Oracle, Nexsan, Fujitsu, Huawei
 - ▶ **Tape:** T10kA, T10kB, T10kC, T10kD, LT08
 - ▶ **No (known) data loss since 2007!**

**Data volume %
(16 PB)**



**Number of files %
(2+ billion)**



Media type	Data volume %	No of files %
Text	28 %	63 %
Images	1 %	7 %
Sound	12 %	10 %
Moving Images	44 %	13 %
Web harvesting	15 %	7 %

Vis alle objekter



Filter

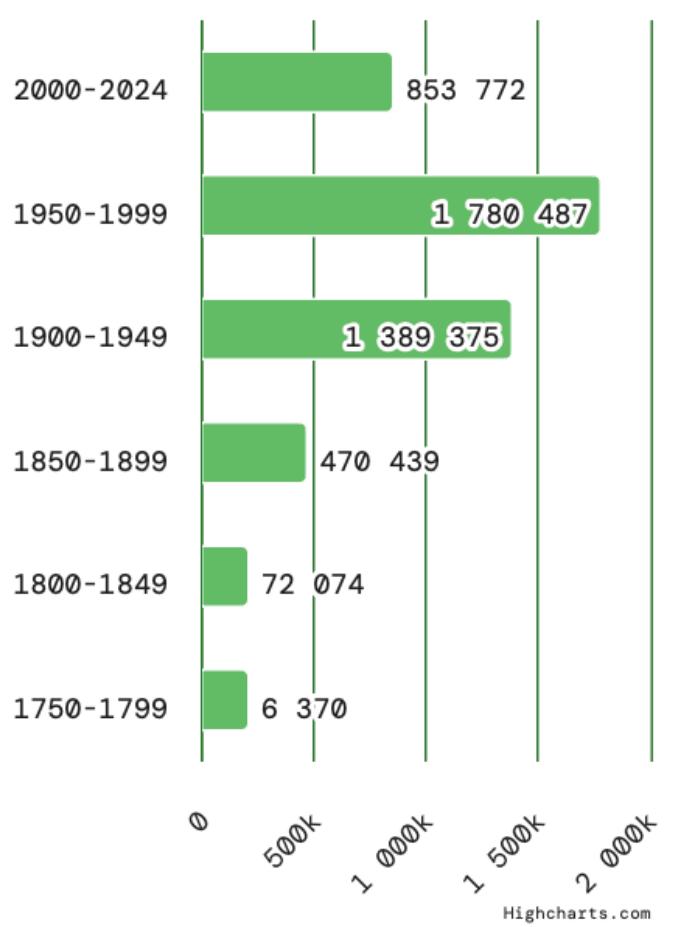
4 572 519 Treff i avis

Relevans ▼ Avisnavn By/fylke

DATO

 Fra dato
dd.mm.åååå Til dato
dd.mm.åååå

Søk



Grimstad
Adressetidende
Torsdag 09.11.1972
Tilgang for alle



VG
Fredag 24.06.1977
Tilgang i norske bibliot...



Lofot-Tidende
Onsdag 02.12.1998
Tilgang for alle



Reform (USA)
Torsdag 07.04.1938
Tilgjengelig etter best...



Folketidende
Tirsdag 23.08.1892
Tilgang for alle



Romerikes Blad
Onsdag 03.04.1968
Tilgang for alle



Firda
Onsdag 15.08.2001
Tilgang for alle



Dalane Tidende
Onsdag 18.08.1965
Tilgang for alle



Kragerø Blad
(Kragerø: 1895-
1997)
Lørdag 02.11.1957
Tilgang for alle



Gjengangeren
Lørdag 21.06.1947
Tilgang for alle



Trønder-fluisa
Torsdag 19.11.2015
Tilgang i norske bibliot...



Bergens tidende
Lørdag 30.11.1963
Tilgang i norske bibliot...



Nye Troms
Torsdag 16.08.1979
Tilgang for alle



Nordlys
Fredag 21.05.1999
Tilgang for alle

<https://nb.no>

Til nb.no

Vis filter

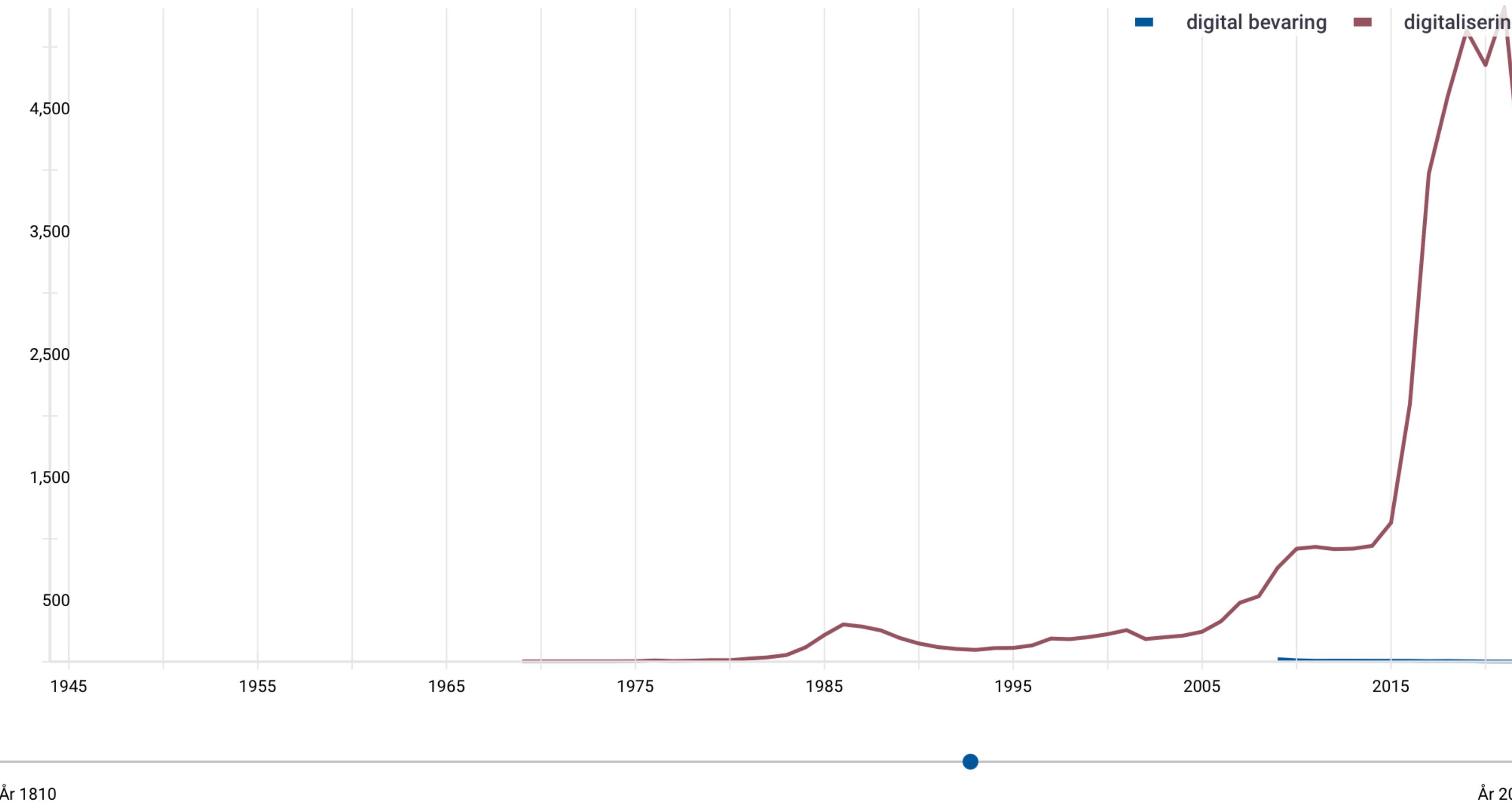
digital bevaring, digitalisering



Korpus: avis, språk: alle målformer

digital bevaring X

digitalisering X



<https://www.nb.no/ngram/>

CHALLENGES: IT

Status →2022

- ▶ Traditional hierarchical organization structure
 - ▶ Decisions floating upwards
- ▶ Myriad stakeholders
 - ▶ Complex prioritizing/competing goals
- ▶ IT bottleneck
 - ▶ Solving problems in an ad hoc manner
 - ▶ Constant context switching
 - ▶ Lack of continuity

Nasjonalbibliotekar - Asia Sira Myhre	
Assisterende nasjonalbibliotekar (Trond Myklebust)	
Digital formidling	Kulturformidling
0104 DICO/DIGR DF - Digital formidling (Trond Myklebust)	0401 KF KF - Kulturformidling (Elne S Kielen)
Biblioteksutvikling	Fag og forskning
0102 BU BU - Biblioteksutvikling - Svein-Arne Timsestrand	0601 FOFO FF - Fag og forskning (Hege Heslein)
Økonomi og personal	Tilvekst og kunnskapsorganisering
0201 ØP ØP - Økonomi og personal (Bernt Fjelberg)	0501 TKL/NTKLO TK - Tilvekst og kunnskapsorganisering (Kjersti Rustad)
0203 OKOO Økonomi (Unni Strauman)	0502 UTVR/UTVO TK utvikling (Hilde Hegdals)
0205 PALO Personal, arkiv, linn (Karl Lie Smith-Meyer)	0506 KUA0 Kunnskapsorganisering monografer (Innre Hole)
Bygg og tekniske tjenester	Kunnskapsorganisering pliktalevrete materialer (Trine Granlund)
0301 BYGO/BYGR BTT - Bygg og tekniske tjenester (Bjørn Skevik)	0507 KUBR DEPOT (Helen Sakrihei)
0303 TEKR/TEKO Teknisk (Solve Bakken)	0511 TRPLUKT TK - Pliktalevring tekst (Malin Prytz)
IT	Tilrettelagt litteratur
0901 IT IT (Wilfred Østgålen)	1101 TL TL Ledelse (Øyvind Engs)
0902 ITP IT Platform (Sverre Bang)	1102 TEKNO Teknologiseksjonen (Arne Kykkjæde)
0903 ITA IT Applikasjon (Frank Sjøvold)	1103 FORM Formellseksjonen (Jørunn Wold)
0904 ITT IT Tjenester (Line Lind)	1104 OUP Offentlig utvalg for punktskrift (Øyvind Engs)
0905 AI AI-løsning (Wilfred Østgålen)	

CHALLENGES: DIGITAL PRESERVATION

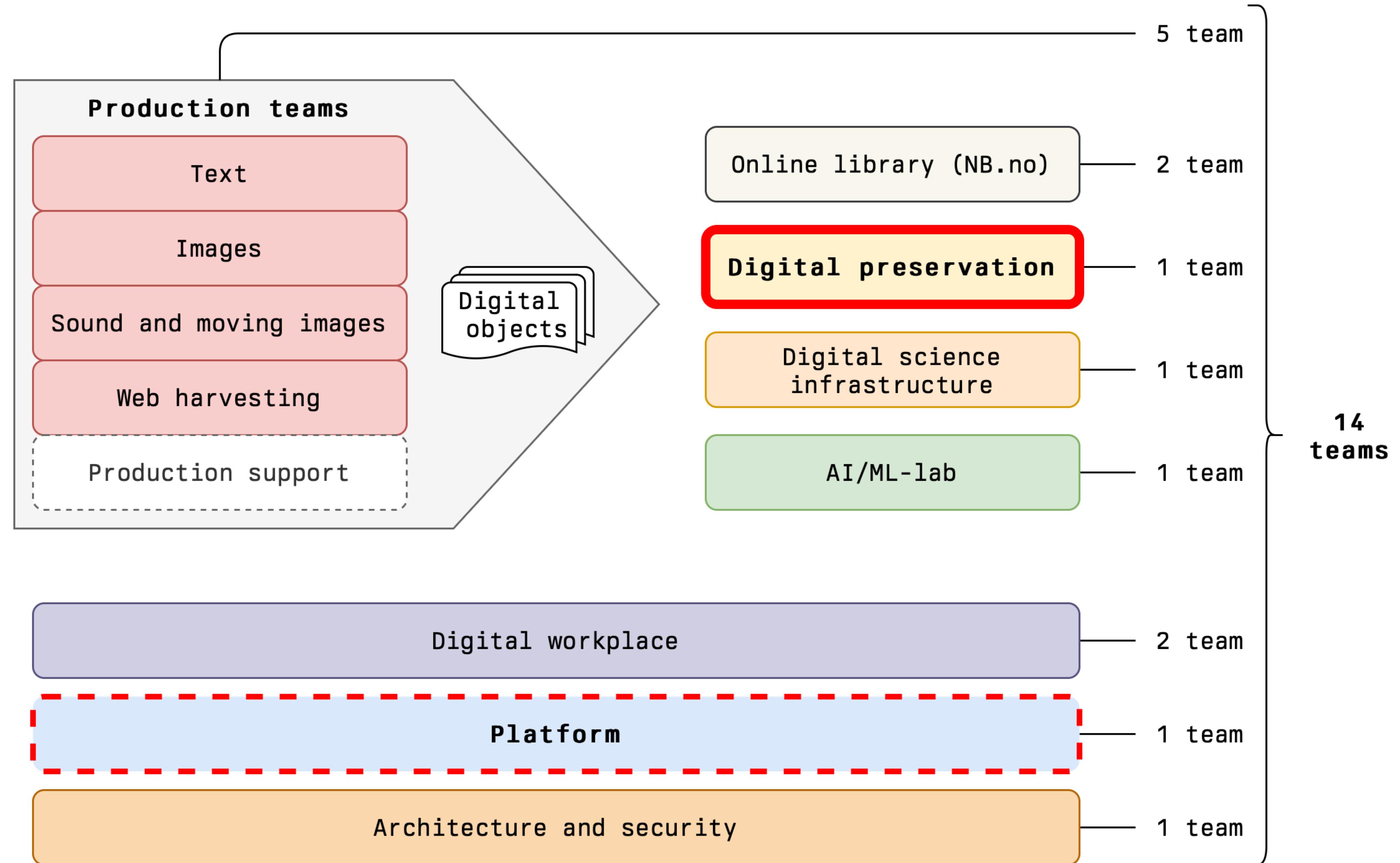
Status →2022

- ▶ Viewed as an IT problem - **previous issues apply**
- ▶ **Store-and-forget** mentality towards files (stored ≠ preserved!)
- ▶ Shared responsibility = **no ones responsibility!**
- ▶ **Lack of consistency:**
 - ▶ Significant amount of files lack of checksums
 - ▶ Lack of metadata
 - ▶ Data packages not standardized
 - ▶ No overall knowledge about what was stored
 - ▶ Limited knowledge about how the data was used

2022 - BIG CHANGE IN IT-STRUCTURE

New IT-director

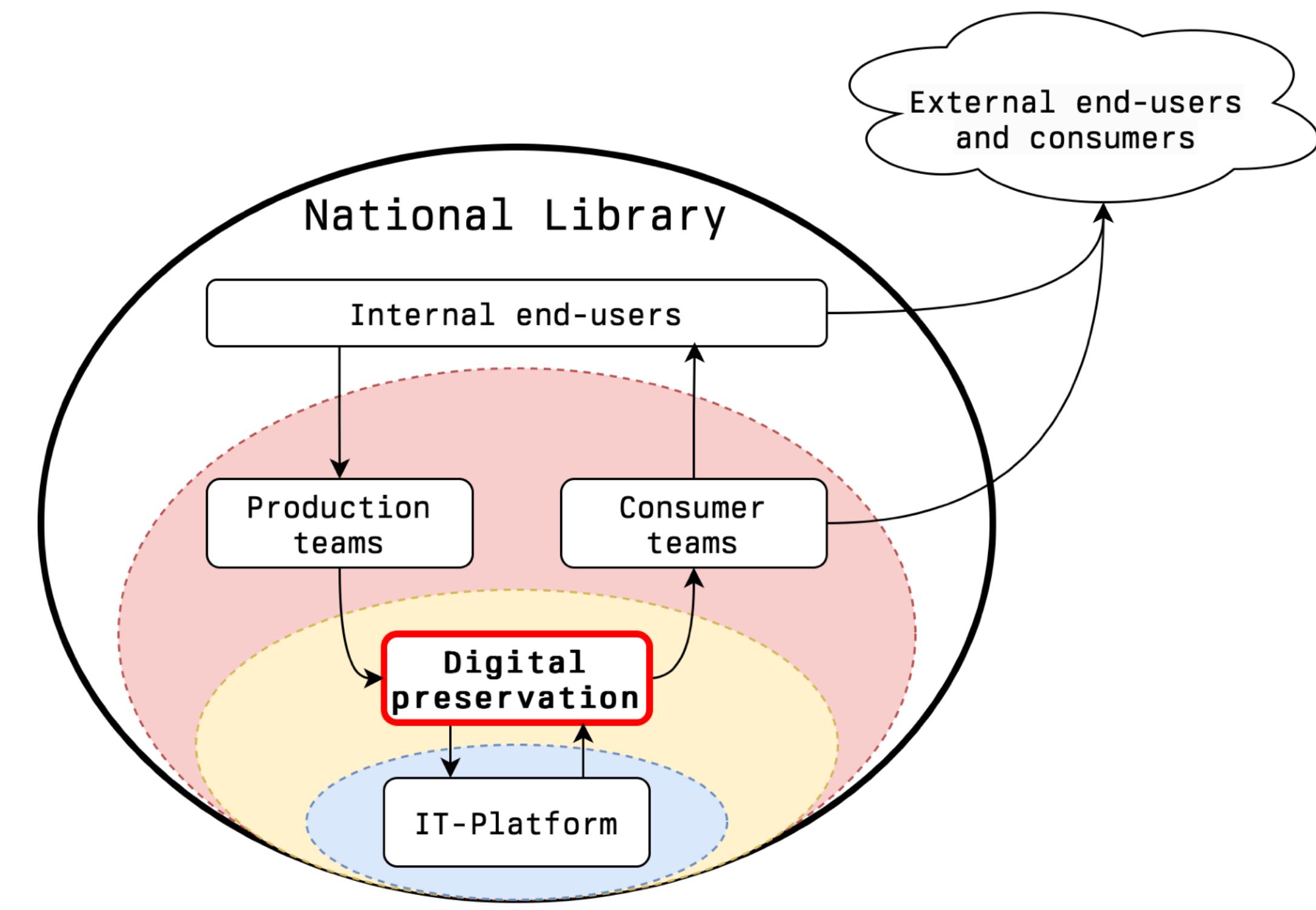
- ▶ **Product orientation** introduced
 - ▶ Build **dynamic products** (MVP) instead of **static solutions**
- ▶ **Autonomous teams** in charge of “products”
 - ▶ **Interdisciplinary** team members (organization looks the same)
 - ▶ Product owner groups **set direction** (directors/section heads)
 - ▶ The team set their own day-to-day **priorities**



DIGITAL PRESERVATION 2022→

Team established June 2022

- ▶ Autonomous, but not interdisciplinary (until late 2023)
 - ▶ Defined scope of responsibility
 - ▶ Where preservation starts (the team does not create data/digitize materials)
 - ▶ Assure files+content lasts (forever)
 - ▶ Primarily an **internal** service
 - ▶ Answers to owner board
 - ▶ Clear priorities
- = Enabling systematic work with digital preservation!



Organization layers/Data flow

DIGITAL PRESERVATION 2022→

As product

Digital preservation team products:

1. Domain expertise of digital preservation in the organization
 - ▶ Responsibility for building competence and spreading awareness
2. “Digital Preservation Services” (DPS)
 - ▶ Responsibility for developing and operating the DPS-software

NLN DIGITAL PRESERVATION STRATEGY

<https://digitalpreservation-blog.nb.no/docs/strategy>

Ambition

- ▶ Ensure the protection of, and meaningful access to, national digital cultural heritage for current and future generations.

Goals for Digital Preservation

- ▶ Digital content for digital preservation shall be received using efficient and standardized machine solutions.
- ▶ Digital content shall be protected against unintended access, alteration, loss, or damage.
- ▶ The National Library shall at all times know what digital content is being preserved, its provenance, its condition, and what has been done to it.
- ▶ Digitally preserved content shall be accessible for dissemination now and in the future.

NLN DIGITAL PRESERVATION PRINCIPLES

<https://digitalpreservation-blog.nb.no/docs/principles>

- ▶ Ensure that digital preservation is done in a sustainable way
- ▶ Use well-documented and open file formats wherever possible
- ▶ Preserve the original file
- ▶ Analyze files that is to be preserved
- ▶ Maintain sufficient metadata to ensure that the files are identifiable and retrievable
- ▶ Use a standardized format to package files for preservation
- ▶ Standardize documentation preservation activities
- ▶ Files should be readable and understandable in the present
- ▶ Ensure that a file is stored in multiple instances, on different storage technologies and in different geographical locations (3-2-1)

Strategy

Home archive policy documents

Roadmap

Ambisjon, mål og strategi for digital bevaring

Published 2024-02-07 · 298 words · Digital Preservation Team | Github source document

Table of Contents

English translation here

Lov om avleveringsplikt for allment tilgjengelege dokument (pliktloven)
§1 [...] vitnemåla om norsk kultur og samfunnsliv kan verta bevarte og gjorde tilgjengelig for alle. Det er viktig at kulturmateriale for forskning og dokumentasjon.¹

Ambisjon for Digital Bevaring:
Sørge for sikring av og meningsfull tilgang til nasjonal digital kulturarv for nåværende og fremtidige generasjoner.

Mål for Digital Bevaring:

- Innhold til digital bevaring skal tas imot med effektive og standardiserte maskinlesebare teknikker.
- Digitalt innhold skal være sikret mot utilstikt utlevering, endring, tap eller skader.
- Nasjonalbiblioteket skal til enhver tid vite hvilket digitalt innhold som blir bevarat, hvilken tilstand det er i og hva som er nørt med det.

Veikart Digital bevaring 2024-2025

Owned by Thomas Edvardsen · Last updated: Feb 09, 2024 by Trond Teigen · 4 min read · 32 people viewed

Strategiske satsingsområder: (T)=teknologi, (S)=standardisering, (K)=Kompetanser

1. Overgang fra SAM-FS til DPS med HPSS som bit-repository (T)

Flytte all daglig tilvekst fra SAM-FS til DPS(HPSS) (30+ lager, 4 TeraByte og 20.000+ kjeldemateriale for forskning og dokumentasjon).
Status finnes i regnearket "Kartlegging av produksjonsløyper" under skillearket "Digital Bevaring".

Fordi:
Det er en nødvendig utfasing av teknisk gjeld. SAM-FS bit-repository har vært leverandøren har meldt EOS (End Of Support) i 2024. Det er innkjøpt og installert en nytt alternativt system (DPS) som erstatter samme funksjonalitet.

2. Flytte historisk materiale fra SAM-FS til DPS med HPSS som bit-repository (T)

Flytte ca.14PetaByte med data fordelt på 16 forskjellige filsystemer i SAM-FS. Det er ikke mulig å flytte alt i én gang, derfor må dataene rearkeiveres i DPS og lagres i HPSS-bit repository. Status finnes i regnearket "Kartlegging av produksjonsløyper" under skillearket "Digital Bevaring".

Fordi:
Det er en nødvendig utfasing av teknisk gjeld. SAM-FS bit-repository har vært leverandøren har meldt EOS (End Of Support) i 2024. Det er innkjøpt og installert en nytt alternativt system (DPS) som erstatter samme funksjonalitet.

Monthly delivery plan

2024 Februar leveringsplan Digital Bevaring

Owned by Trond Teigen · Last updated: Feb 09, 2024 · 1 min read · 5 people viewed

- Produksjonsløype for re-arkivering av DSM-materiale (Radio) fra SAM-FS til DPS. Genererer MP4-visningsformat og bruk av eArk som pakkeformat. **Veikart:** Punkt 2 Flytte historisk materiale fra SAM-FS til DPS
- Rearkivering av aviser fra SAM-FS, disse mangler sjekksummer som er nødvendig for framtidig bruk. **Veikart:** Punkt 2 Flytte historisk materiale fra SAM-FS til DPS
- Beskrive scenarier for "eierflagging" av samlinger og hva som skal skje med dem. Danner grunnlag for implementering av "eierskap" og tilgangskontroll. **Veikart:** Punkt 3.G Eierskap og tilganger til bevart materiale
- Lage oversikt over hvilke typer metadata Digital bevaring har behov for. Danner grunnlag for framtidig bruk. **Veikart:** Punkt 3.A Definere metadataformat for SIP og punkt 3.I gjenomføring
- Gjennomføre ROS analyse for Digital bevaring. Danner grunnlag for opprettelse av DPS. **Veikart:** Punkt 3.B Flytte historisk materiale fra SAM-FS til DPS

Main activities (epics)

- DB-645 Techtalk: eArchiving (e-ark) som standardformat for bevaring
- DB-651 Drift av DPS
- DB-667 Finne ut hvordan vi kan motta flere filtre i utelevering, slik at vi kan filtrere etter dato, type osv.
- DB-485 Endre parameter context til å bruke arv
- DB-650 Test om Siegfried håndterer store filer bedre enn Droid
- DB-682 Kan vi erstatte DroidIdentificationProcessor med Siegfried
- DB-666 Finne ut hvorfor Droid ikke identifiserer wav-filer riktig
- DB-683 Oppdatere signaturfil for Droid
- DB-613 Sende DPS logger til Logstash/ES/Kibana
- DB-585 NB-pakker til E-ARK
- DB-587 Rearch DSM Radio
- DB-607 Lage Pronom signatur for JSON Lines
- DB-624 Rearkivering av historisk radiomateriale i DPS
- DB-622 Teste rearkeivering flyt ende til ende
- DB-662 Finne ut om vi kan lage MD5 sjekksum for Representasjoner
- DB-665 Bytte ut Droid med Siegfried
- DB-688 MP3 visningsfiler skal ikke flyttes fra DSM til DPS

Kanban board (tasks)

Digital Bevaring

TIL UTFØRING 12 · UNDER ARBEID 10 · KLAR TIL VERIFISERING 10

OCFL: Felge opp forslag fra Jürgen Enge · REARK DSM RADIO · DB-622

DIVERSE · DB-556

Endre parameter context til å bruke arv · DRIFT AV DPS · DB-485

DRIFT AV DPS · DB-689

Dokumentere hvilke verktøy som er aktuelle å bruke for validering og metadata extraction · DIVERSE · DB-454

Gjøre tester av rawcooked for encoding av DPX-sekvenser · DB-689

Blogpost Strategi · STRATEGI DIGBEV · DB-685

Test om Siegfried håndterer store filer bedre enn Droid · DRIFT AV DPS · DB-650

Validere E-ARK sip · REARK DSM RADIO · DB-650

BUILDING DOMAIN EXPERTISE

- ▶ Share experiences and policy documents on our [blog](#)
- ▶ Revision [preferred file format list](#)
- ▶ Membership **Digital Preservation Coalition (DPC)**
 - ▶ DPC assessment tools (DPC-RAM, DPC-CAT)
 - ▶ DPC bitlist council
- ▶ Involvement in national and international community

The screenshot shows a blog post titled "Digital preservation at the National Library of Norway". The header includes links for Home, archive, policy documents, search, about, and nb.no. The main content area displays the blog post with sections for Strategy, Roadmap, Monthly delivery plan, Main activities (epics), and Kanban board (tasks). Below the post is a summary: "Ambitions, Goals, and Strategy for Digital Preservation at the National Library". It states that the Digital Preservation Team at the National Library (NLN) has developed its first strategy for digital preservation. This strategy aims to steer, structure, and sharpen the... Published 2024-02-20 · 3 min · 626 words · Trond Teigen.

NiFi S2S on Secured Instances
Guide to setting up a Site-to-Site (S2S) communication between two secured NiFi instances with user and policy management. This guide is based on experiences from... Published 2024-02-16 · 13 min · 2654 words · Daniel Aaron Salwerowicz

BIT REPOSITORY REPLACEMENT

2020-2022

- ▶ SAM-FS EOL (2021)
- ▶ Tender (2021) → IBM High Performance Storage System (HPSS)
- ▶ Installed (2022)

CLOUD VS. IN-HOUSE?

- ▶ Cons of cloud:
 - ▶ Possible **performance** challenges when moving large data volumes
 - ▶ **Legal uncertainties** in relation storing cultural heritage materials at commercial vendors, potentially outside of Norwegian borders
 - ▶ The costs of **retrieving** large amount of data from cloud provider
 - ▶ **Lack of in-house experience** with cloud infrastructure at the time
 - ▶ Solid **in-house experience with self-hosting** (we were comfortable to keep on doing it)

HPSS (IBM)

High Performance Storage System

- ▶ Linux OS
- ▶ Block based storage
- ▶ Mix of different disk and tape technologies (no vendor lock-in)
- ▶ **Disk:** Fujitsu, Huawei og Nexsan
- ▶ **Tape:** LT08 in 2 SL8500 (10k slot libraries)



SELECTION OF HPSS

High Performance Storage System

- ▶ Supports 3-2-1 (disk+tape+tape)
- ▶ Scales well
- ▶ No vendor lock-in (multi-vendor HW)
- ▶ Multilevel checksumming (blocks and files)
- ▶ Large user community
 - ▶ 30+ clients and 3+Exabytes stored in HPSS systems worldwide



HPSS IMPLEMENTATION

Window of opportunity

- ▶ Installation → Production lines still writing to SAM-FS
- ▶ Move bits as-is from SAM-FS to HPSS?
- ▶ Opportunity to do things better!
- ▶ Establish new ingest/preservation/dissemination methodology according to principles

OFF-THE-SHELF PRODUCTS?

2021-2022

- ▶ **Criteria:**

- ▶ Handles **large data volume** and expected growth
- ▶ Need for **automated processes**
- ▶ Has **separated preservation and playback** in different solutions
- ▶ **Standardized and open solutions**
- ▶ Surveyed the market (looked at Archivematica, Libnova, CSC, and more)
- ▶ **None of these fit our needs – challenges with:**
 - ▶ Scale in **data volume**
 - ▶ **Licensing** (often volume-based)
 - ▶ **Running environment** that does not fit in NLN architecture
 - ▶ Systems contained **functionality** that NLN did not request (viewing/playback)

DIY → DPS

“Digital Preservation Services”

- ▶ Developed by digital preservation team (Jun 2022 → Dec 2022)
- ▶ Built after preservation principles
- ▶ DPS 1.0 = Unified ingest workflow to HPSS
 - ▶ Checksums for all files (and SIPs)
 - ▶ Stored in HPSS (along with files), in HPSS (DB2 database), and in DPS (locationDB)
 - ▶ Standardized delivery format (not package format)
 - ▶ Asynchronous communication
 - ▶ Inventory database with information on:
 - ▶ Content type, File types, Number of files in package, Location in bit repository, Events regarding ingest, Who delivered SIP

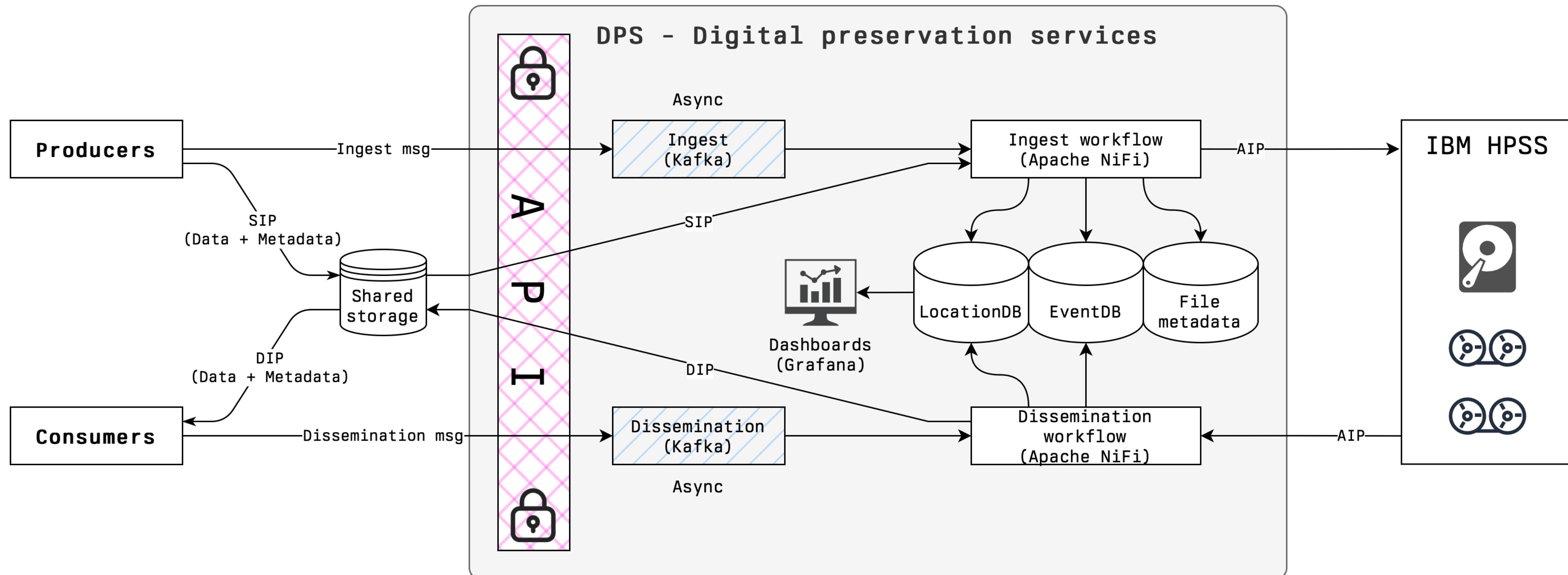
DPS DEVELOPMENT

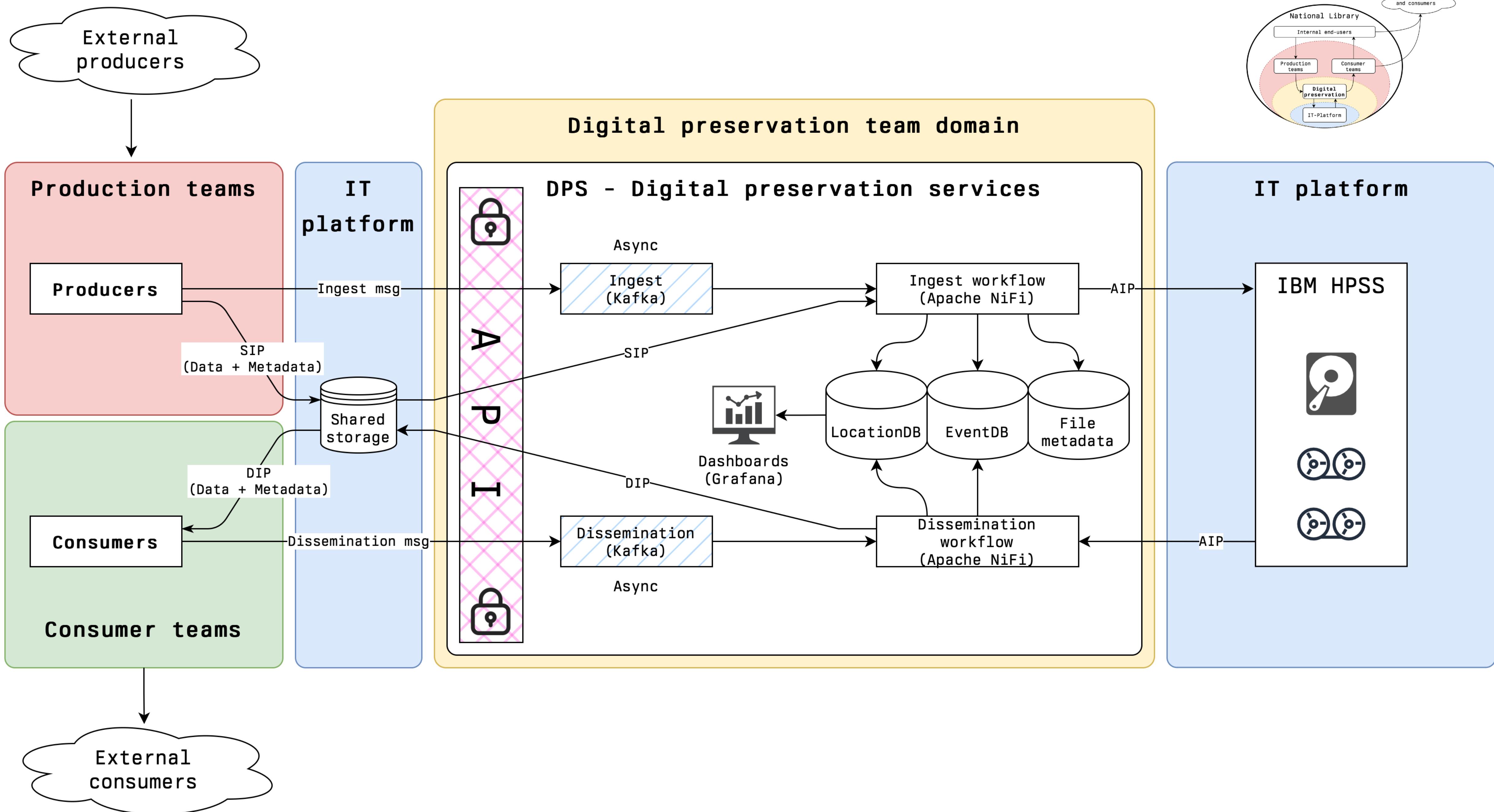
2022→2024

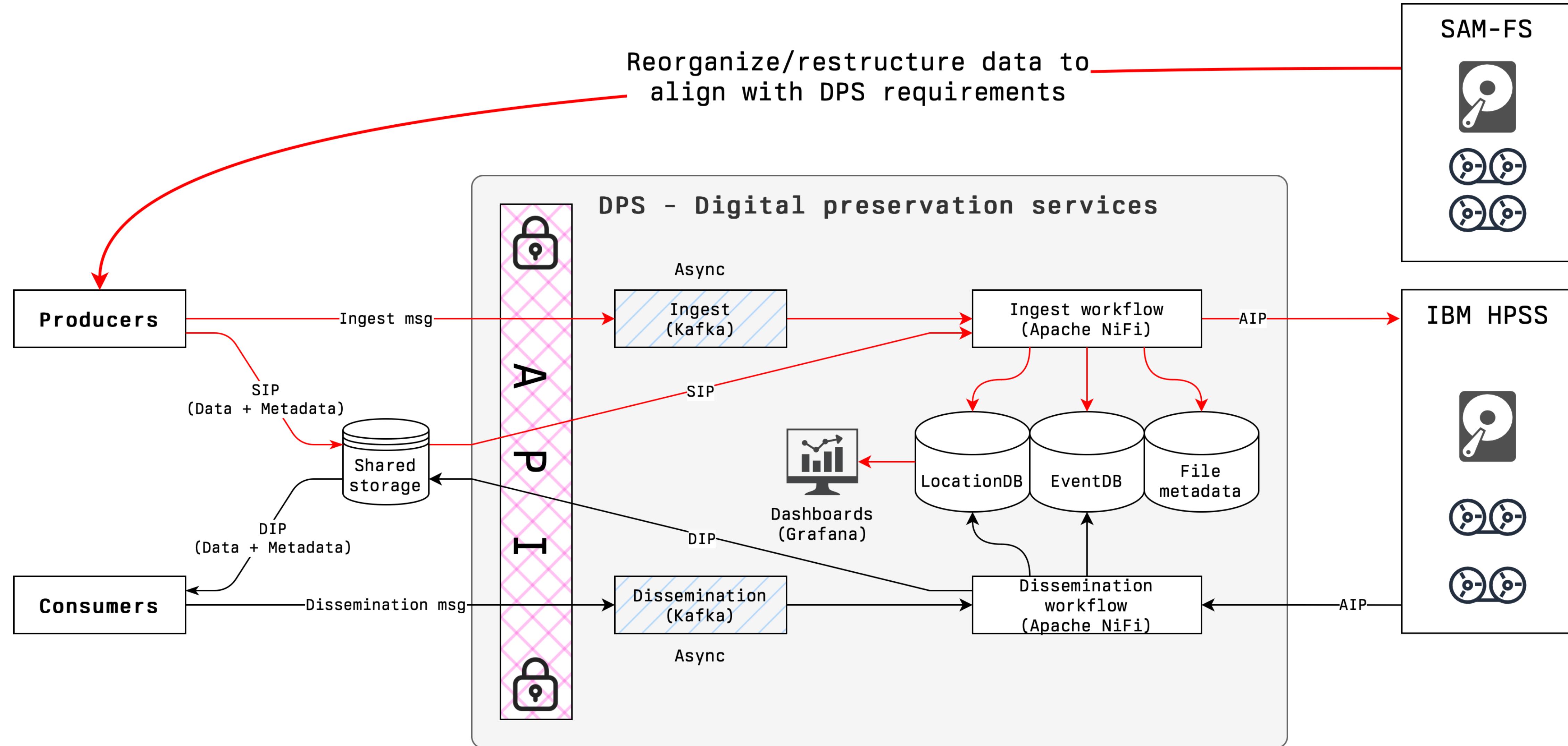
- ▶ **Iterative** and **incremental** development (MVP):
 - ▶ Addition of dissemination workflow (2023)
 - ▶ Ingest workflow expanded (2023→)
 - ▶ File identification (DROID/Siegfried)
 - ▶ File validation etc.
- ▶ Update 30+ production lines to deliver new data to DPS instead of SAM-FS.
 - ▶ Took most of 2023 to accomplish (1 production line remaining still!)

DPS TECHNOLOGIES

- ▶ **Java, Spring Boot, Keycloak, Kubernetes** for REST APIs for SIP+DIP messages
- ▶ **NFS/GlusterFS** shared storage for transferring SIP+DIP packages
- ▶ **Apache Kafka** for asynchronous transfer of messages
- ▶ **Apache NiFi** for processing SIP/AIP/DIP packages
- ▶ **IBM High Performance Storage System (HPSS)** for archival storage
- ▶ **Grafana** dashboards for monitoring and statistics

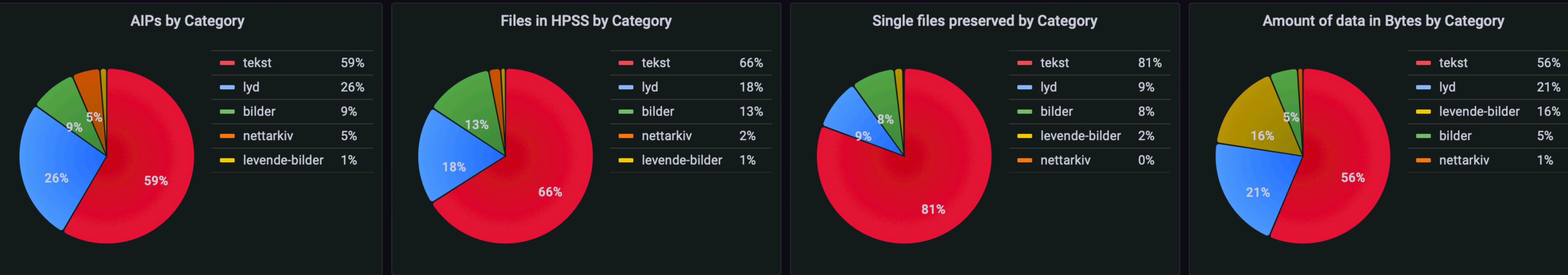






REARCHIVING SAM-FS → HPSS (2023→2025)

Number of AIPs and Files, and amount of Data by Category and Type



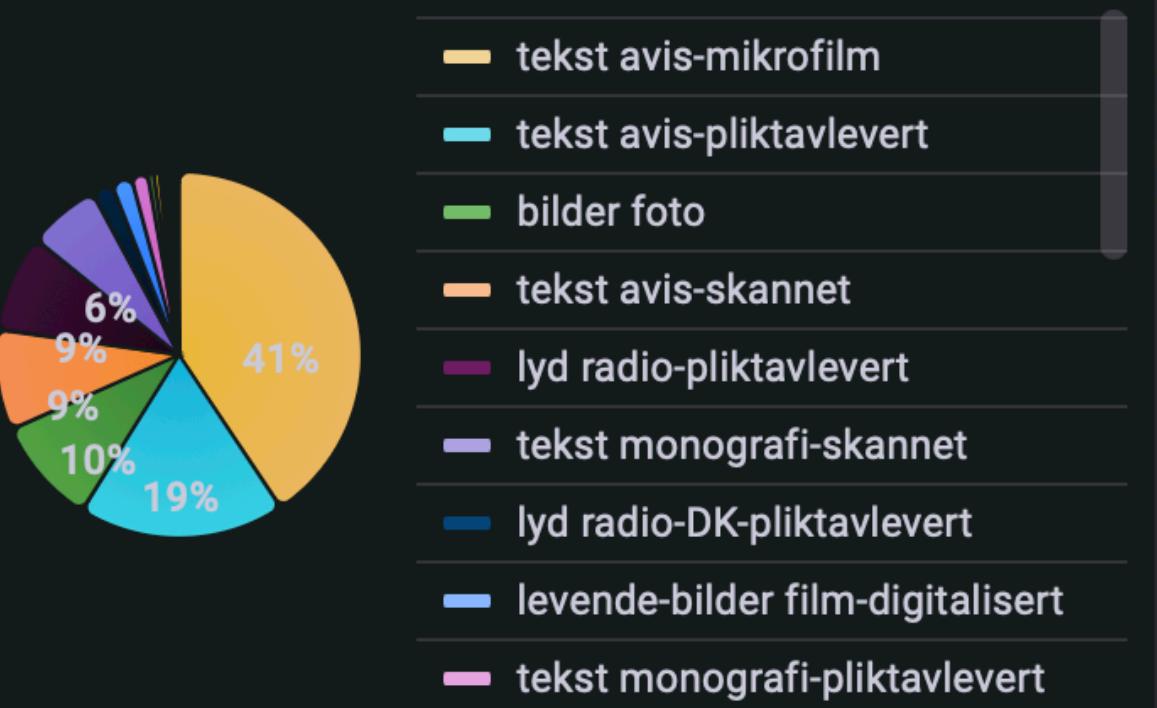
AIPs by Category and Type



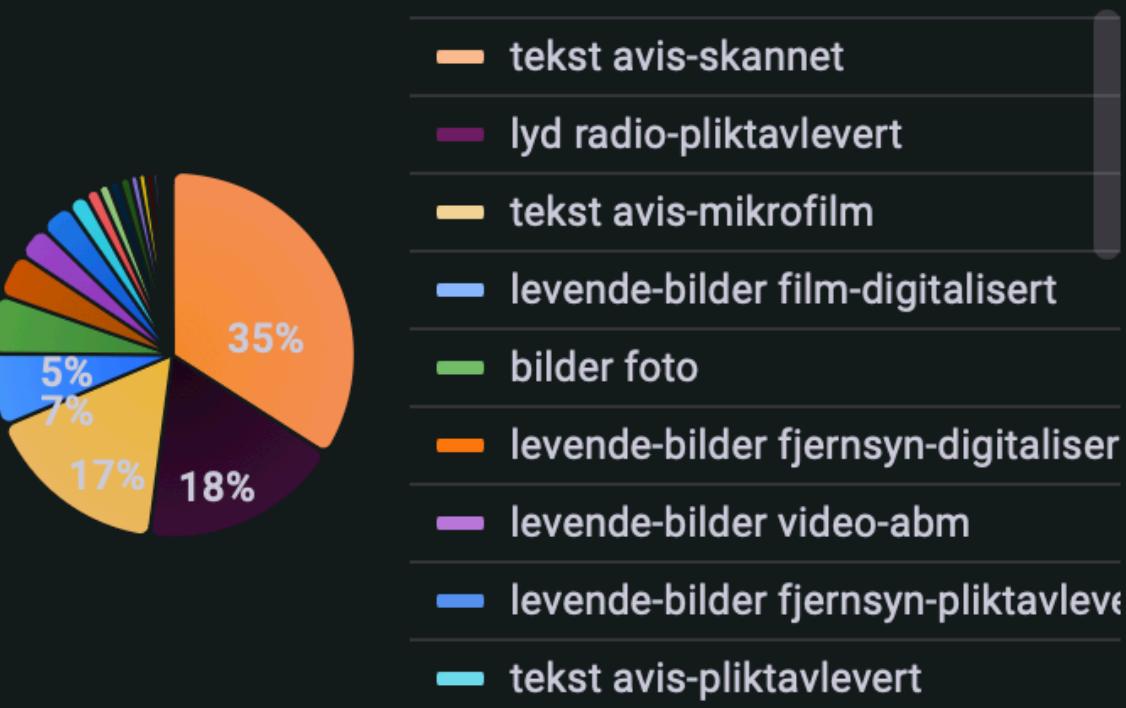
Files in HPSS by Category and Type



Single files preserved by Category and Type



Amount of data in Bytes by Category and Type



AIPs by Category and Type

Category	Type	Count ↓
tekst	avis-skannet	2,382,151
tekst	avis-mikrofilm	1,898,377
lyd	radio-pliktavlevert	1,686,904
bilder	foto	711,289
nettarkiv	warc	424,441
lyd	radio-DK-pliktavlevert	225,777

Files in HPSS by Category and Type

Category	Type	Count ↑
nettarkiv	acquisition	21
levende-bilder	videokunstarki...	1,260
lyd	musikk-studio-f...	4,260
levende-bilder	film-digitalisert	5,227
levende-bilder	fjernsyn-DK-pli...	6,816
levende-bilder	fjernsyn-digitali...	8,204

Single files preserved by Category and Type

Category	Type	Count ↓
tekst	avis-mikrofilm	104,184,196
tekst	tidsskrift-skannet	47,211,872
tekst	avis-pliktavlevert	46,641,050
bilder	foto	24,115,526
tekst	avis-skannet	22,608,577
lyd	radio-pliktavlevert	21,842,257

Amount of data in Bytes by Category and Type

Category	Type	Bytes ↓
tekst	avis-skannet	1.50 PiB
lyd	radio-pliktavlevert	804.87 TiB
tekst	avis-mikrofilm	745.98 TiB
levende-bilder	film-digitalisert	288.14 TiB
bilder	foto	233.07 TiB
levende-bilder	fjernsyn-digitalis...	176.12 TiB

▼ Total number of AIPs and Files, and amount of Data preserved

Total number of AIPs preserved

8,210,715

i Total number of files stored in HPSS

39,450,431

i Total single files preserved

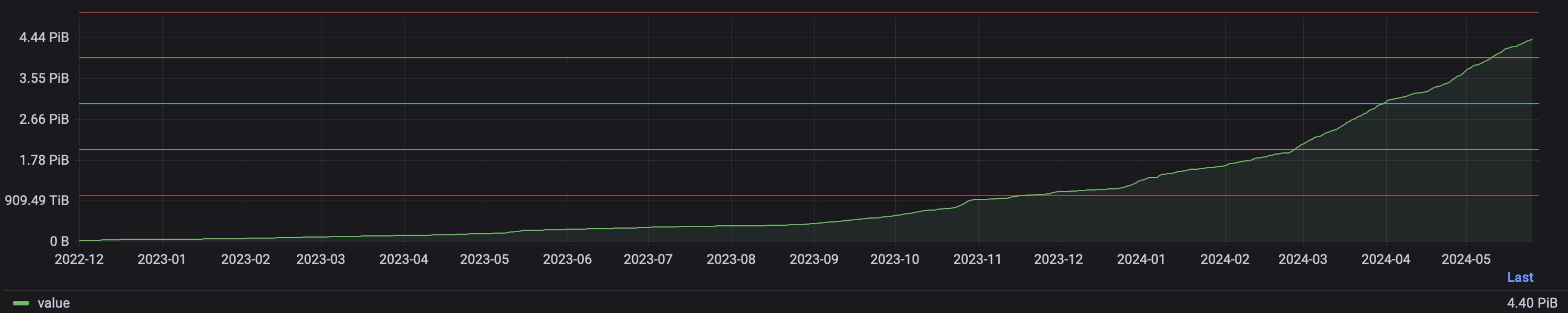
298,141,949

Total amount of data preserved

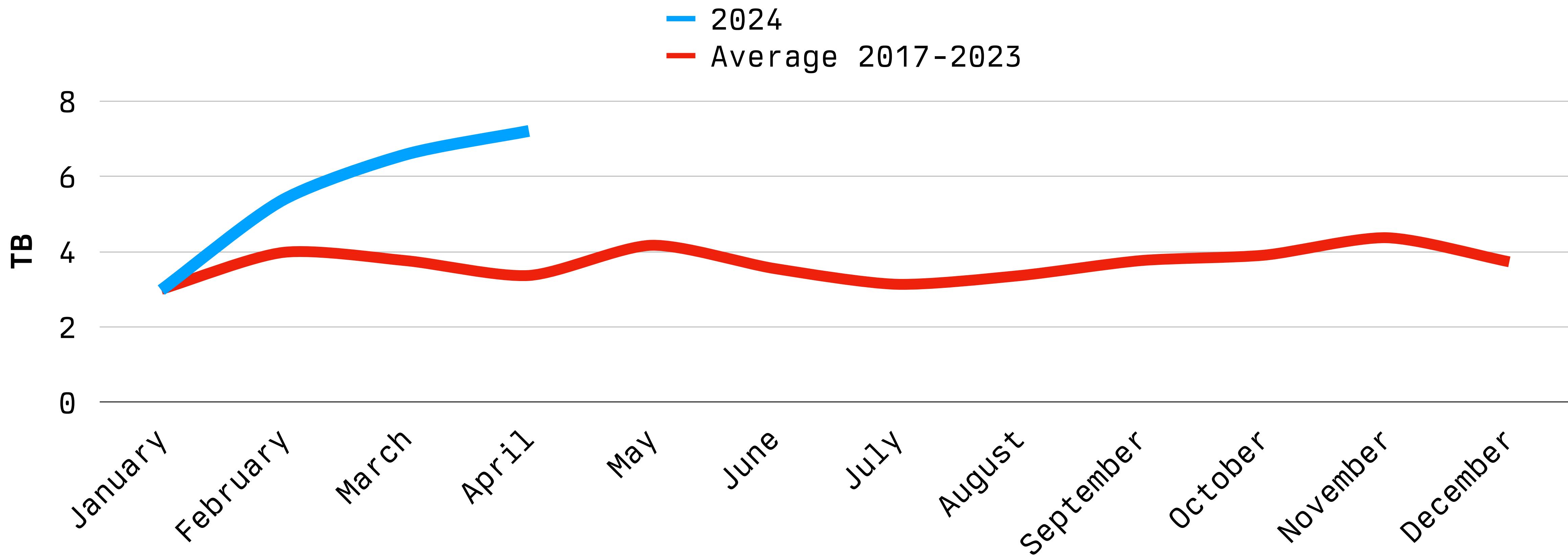
4.42 PiB

▼ Total bytes preserved in DPS

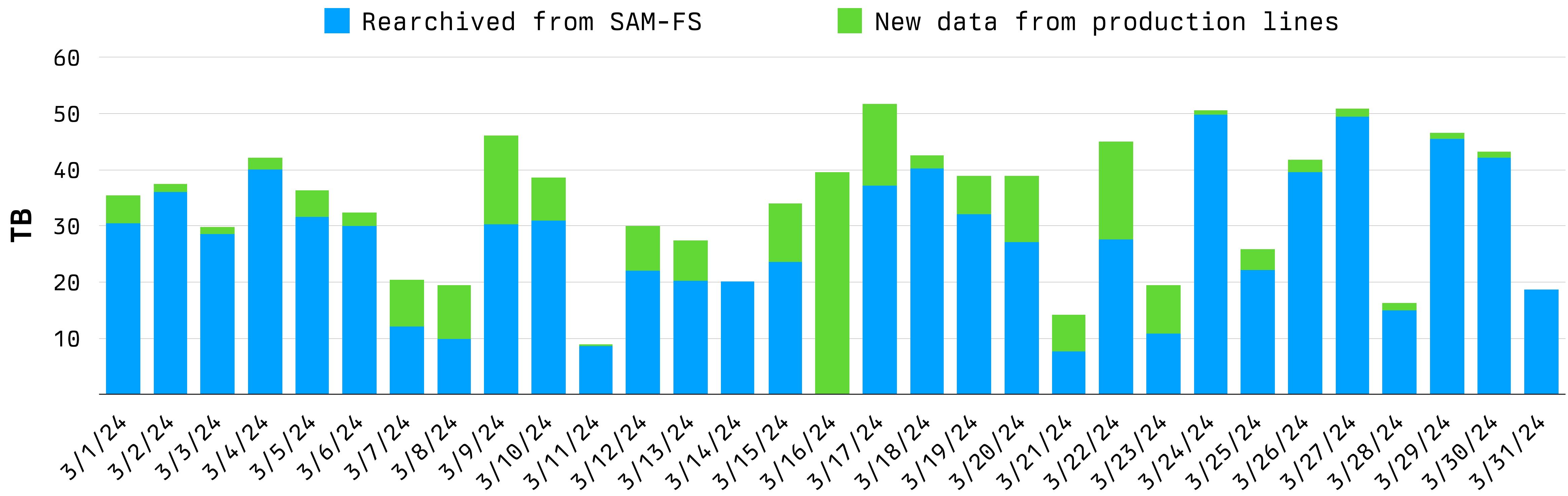
Total Bytes Preserved in DPS



AVERAGE DAILY INGEST OF NEW DATA PER MONTH IN TB



DAILY INGEST TO DPS IN TB, MARCH 2024



Ingest in TB In March	Total Re-archived from SAM-FS	New data
Accumulated	1 035	838
Daily average	33	27

DPS 2.0 - FUTURE PLANS

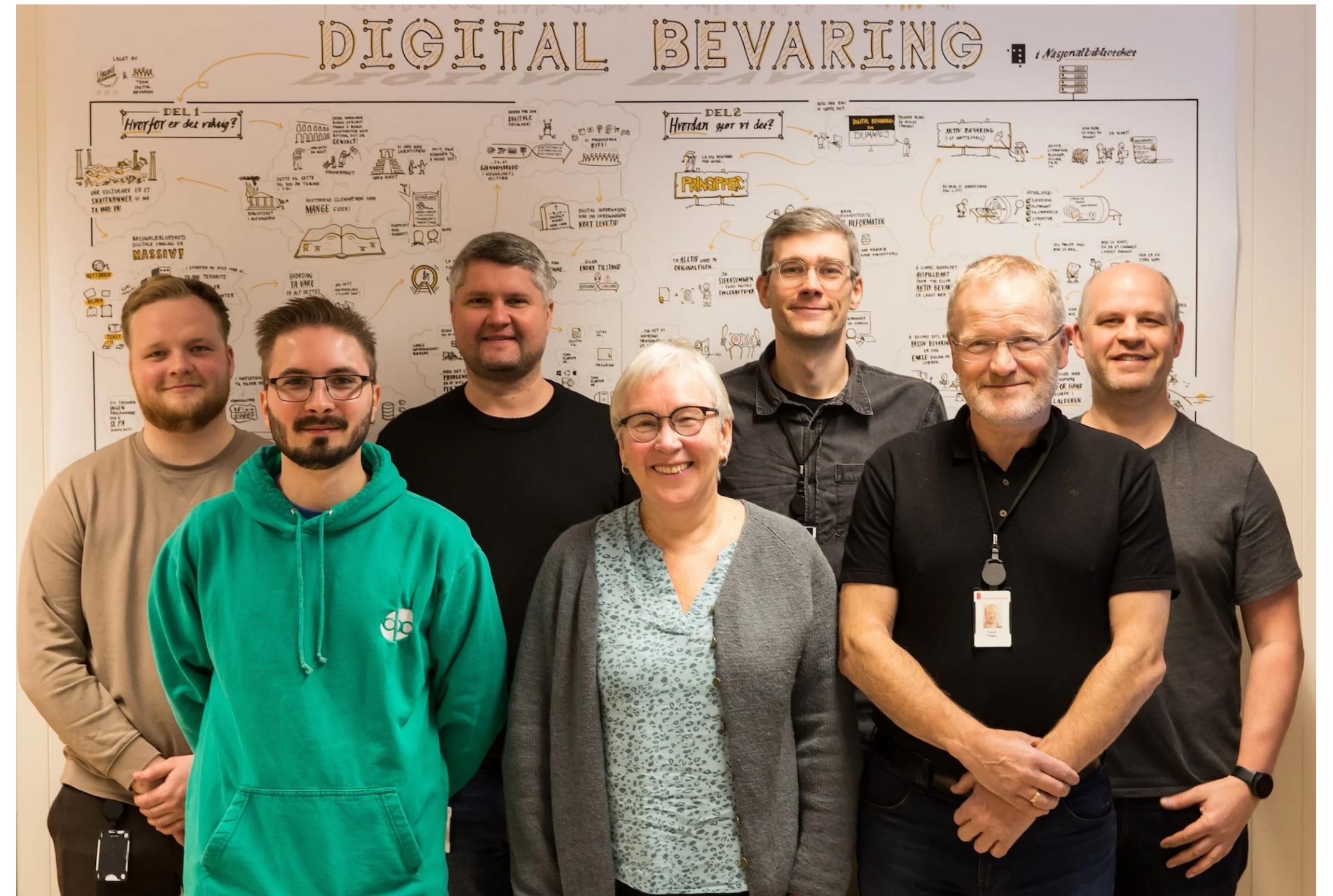
2024→2025

- ▶ Keep improving in small increments!
- ▶ Authentication and Authorization
 - ▶ Lock down the DPS (WIP)
 - ▶ Role-based access (near future)
- ▶ Standardize information package content structure
 - ▶ Implement eArchiving Standards & Specifications (WIP)
- ▶ Improve ingest workflow to handle unpacked files
 - ▶ Get control at file level (not .tar level)



CONTACT INFO

- ▶ torbjorn.pedersen@nb.no
- ▶ digitalpreservation-blog.nb.no/
- ▶ NB.no/



The digital preservation team