# HDF5: The most versatile container for sharing scientific and engineering data

Gerd Heber, The HDF Group

# Key Points

A big **thank you** to DESY, the Organizing Committee, our sponsors, and countless helpers

- HDF5 is well-positioned as a versatile and evolving container for structured scientific and engineering data
- Like street lights, HDF5 is a <u>public good</u> (non-rivalrous, non-excludable)
- It's evolution is intertwined with another public good: <u>language</u>
- These are exciting times for anyone interested in either or both!
- Public goods face unique challenges: Let's work them together!

Facts inform. Logic convinces. But stories move us. So let me tell you a story.
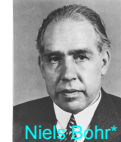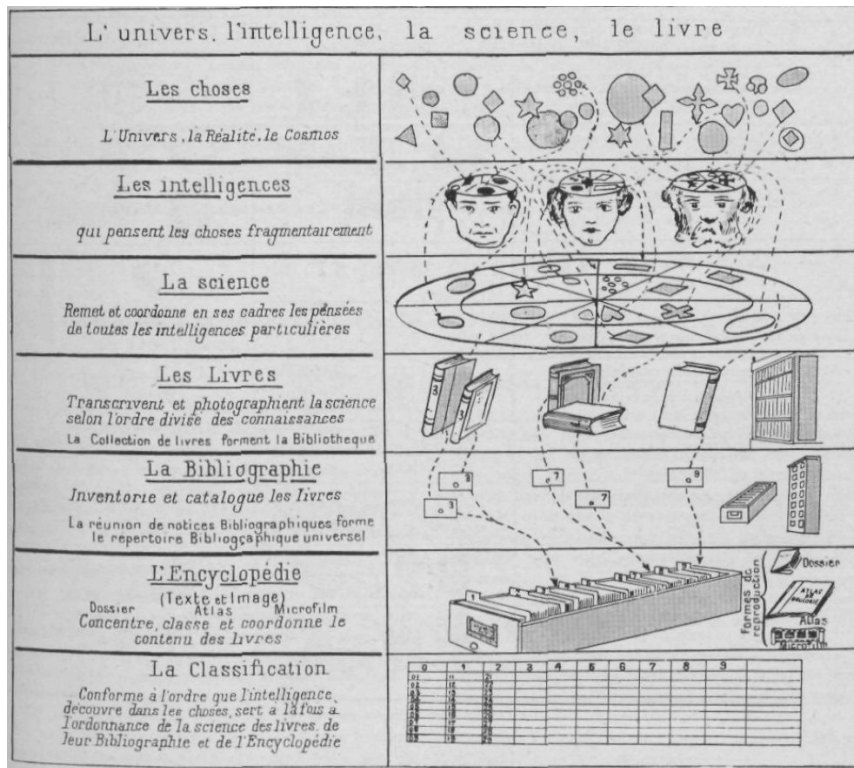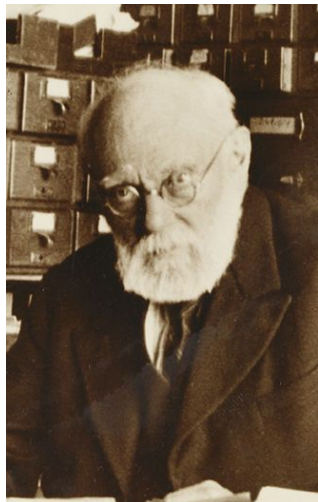
# What could possibly go wrong?



Vienna, 1900

"license with dates and details ... has been, though not unsparingly, indulged" (Edward Bulwer-Lytton)

Paul Otlet
*Traité de documentation: le livre sur le livre, théorie et pratique*
Brussels, 1934



L' univers, l'intelligence, la science, le livre

Les choses
L'Univers, la Réalité, le Cosmos

Les intelligences
qui pensent les choses fragmentairement

La science
Remet et coordonne en ses cadres les pensées
de toutes les intelligences particulières

Les Livres
Transcrivent et photographient la science
selon l'ordre divisé des connaissances
La Collection de livres forment la Bibliothèque

La Bibliographie
Inventorie et catalogue les livres
La réunion de notices Bibliographiques forme
le répertoire Bibliographique universel

L'Encyclopédie
(Texte et Image)
Dossier        Atlas        Microfilm
Concentre, classe et coordonne le
contenu des livres

La Classification
Conforme à l'ordre que l'intelligence
découvre dans les choses, sert à la fois à
l'ordonnance de la science des livres, de
leur Bibliographie et de l'Encyclopédie

* NOT true!

Niels Bohr*

Ludwig Wittgenstein     Adolf Hitler

?

Vienna, 1900

Flakturm IV, Heiligengeistfeld, Hamburg, 1942
(~ 8 km southeast of DESY)

* NOT true!

Niels Bohr*

Ludwig Wittgenstein

Adolf Hitler

?

Vienna, 1900

Trinity, 1945
Ivy Mike, 1952

H-BOMB VS A-BOMB:
WHICH IS MORE POWERFUL?

L'univers, l'intelligence, la science, le livre

Les choses
L'Univers, la Réalité, le Cosmos

Les intelligences
qui pensent les choses fragmentairement

La science
Remet et coordonne en ses cadres les pensées
de toutes les intelligences particulières

Les Livres
Transcrivent et photographient la science
selon l'ordre divisé des connaissances
La Collection de livres forment la Bibliothèque

La Bibliographie
Inventorie et catalogue les livres
La réunion de notices Bibliographiques forme
le répertoire Bibliographique universel

L'Encyclopédie
(Texte et Image)
Dossier      Atlas      Microfilm
Concentre, classe et coordonne le
contenu des livres

La Classification
Conforme à l'ordre que l'intelligence
découvre dans les choses, sert à la fois à
l'ordonnance de la science des livres, de
leur Bibliographie et de l'Encyclopédie

# HDF5 1.0.0 November 6, 1998

* NOT true!

Niels Bohr*

Ludwig Wittgenstein

Adolf Hitler

?

Vienna, 1900

L' univers, l'intelligence, la science, le livre

Les choses

L'Univers, la Réalité, le Cosmos

Les intelligences

qui pensent les choses fragmentairement

La science

Remet et coordonne en ses cadres les pensées de toutes les intelligences particulières

Les Livres

Transcrit et photographient la science selon l'ordre des connaissances. La Collection de livres forme la Bibliothèque.

La Bibliographie

Inventorie et catalogue les livres. La réunion des tables Bibliographiques forme la table Bibliographique universel.

L'Encyclopédie
(Texte et Image)
Dossier    Atlas    Microfilm
Concentre, classe et coordonne le contenu des livres.

La Classification.

Conforme à l'ordre que l'intelligence découvre dans les choses, sert à la fois à l'ordonnance de la science des livres, de leur Bibliographie et de l'Encyclopédie.

Element Types    Basis Functions and Interpolation Schemes    sparse and dense fields    Field value types

Mesh Types    Coordinate Systems    Storage Conventions And Data Structures    Fields of Fields of Fields

Mesh Decompositions    Compression

HDF5

HDF5    group    group    datasets
metadata
metadata    metadata    metadata
metadata
group    dataset
metadata    metadata
metadata

https://www.youtube.com/@hdf5

# Wittgenstein's Turn and Quantum 2nd



L' univers, l'intelligence, la science, le livre

**Les choses**
L'Univers, la Réalité, le Cosmos

**Les intelligences**
qui pensent les choses fragmentairement

**La science**
Remet et coordonne en ses cadres les pensées de toutes les intelligences particulières

**Les Livres**
Transcrit et photographient la science selon l'ordre des connaissances. La Collection de livres forme la Bibliothèque.

**La Bibliographie**
Inventorie et catalogue les livres. La réunion des notices Bibliographiques forme le répertoire Bibliographique universel

**L'Encyclopédie** (Texte et Image)
Dossier · Atlas · Microfilm
Concentre, classe et coordonne le contenu des livres

**La Classification**
Conforme à l'ordre que l'intelligence découvre dans les choses, sert à la fois à l'ordonnance de la science des livres, de leur Bibliographie et de l'Encyclopédie

* NOT true!

Niels Bohr*

Ludwig Wittgenstein · Adolf Hitler

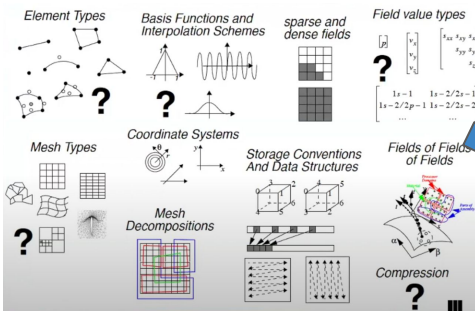$P(w_n|w_1, w_2, \ldots, w_{n-1})$ *

"People's language reflects their reality;
LLM's language reflects people's language."

* What's next?

$W_{pot}$  $\Delta x$

$|\Psi|^2$

$-a$  $0$  $a$  $x$

HDF5

Element Types · Basis Functions and Interpolation Schemes · sparse and dense fields · Field value types

Mesh Types · Coordinate Systems · Storage Conventions And Data Structures · Fields of Fields of Fields

Mesh Decompositions · Compression

group · group · datasets · metadata

metadata · metadata · metadata

HDF5

group · dataset · metadata

metadata · metadata
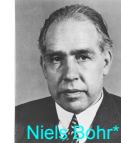
"It is the stillest voice that brings on the storm. Thoughts that come on doves' feet guide the world." (Friedrich Nietzsche)

# X-ready HDF5?



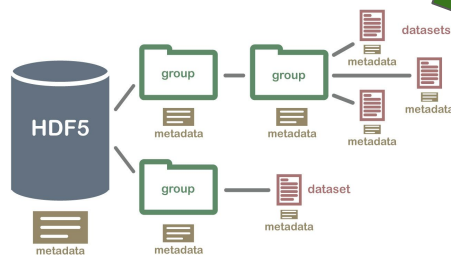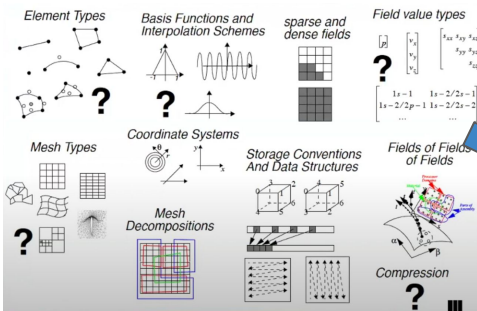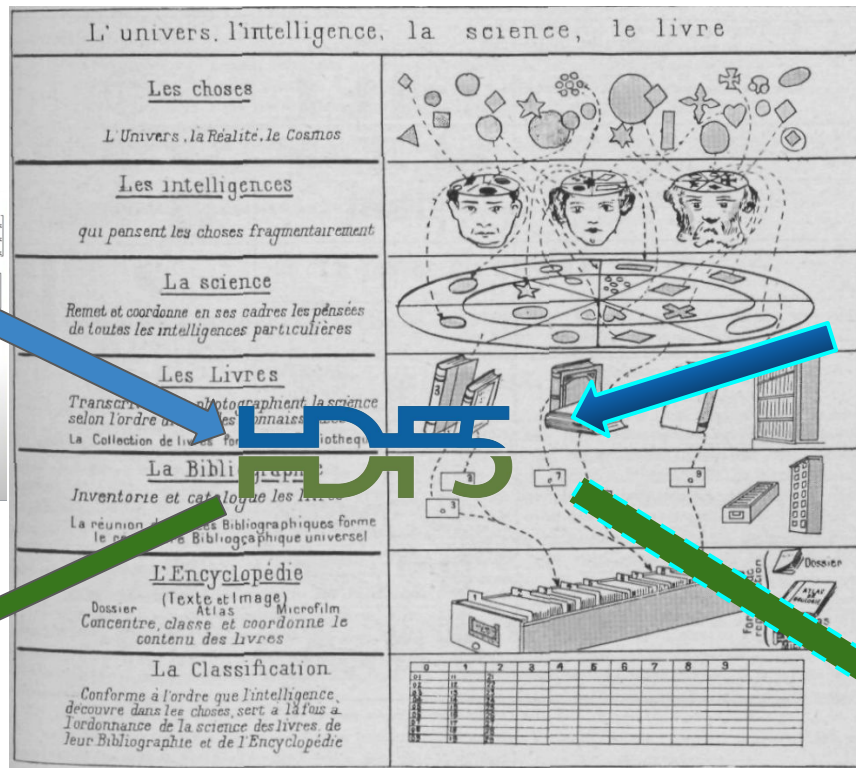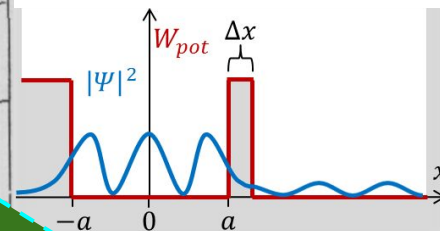* NOT true!

Niels Bohr*

Ludwig Wittgenstein

Adolf Hitler

$P(w_n|w_1, w_2, \ldots, w_{n-1})$ *

"People's language reflects their reality;
LLM's language reflects people's language."

* What's next?

HDF5

Element Types

Basis Functions and Interpolation Schemes

sparse and dense fields

Field value types

Mesh Types

Coordinate Systems

Storage Conventions And Data Structures

Fields of Fields of Fields

Mesh Decompositions

Compression

L'univers, l'intelligence, la science, le livre

Les choses
L'Univers, la Réalité, le Cosmos

Les intelligences
qui pensent les choses fragmentairement

La science
Remet et coordonne en ses cadres les pensées de toutes les intelligences particulières

Les Livres
Transcrit et photographie la science selon l'ordre des connaissances
La Collection de livres forme la Bibliothèque

La Bibliographie
Inventorie et catalogue les livres
La réunion des Notices Bibliographiques forme la Bibliographie universel

L'Encyclopédie
(Texte et Image)
Dossier    Atlas    Microfilm
Concentre, classe et coordonne le contenu des livres

La Classification
Conforme à l'ordre que l'intelligence découvre dans les choses, sert à la fois à l'ordonnance de la science des livres, de leur Bibliographie et de l'Encyclopédie

$|\Psi|^2$    $W_{pot}$    $\Delta x$

$-a$    0    $a$    $x$

HDF5

group

group

datasets
metadata

metadata

metadata

metadata

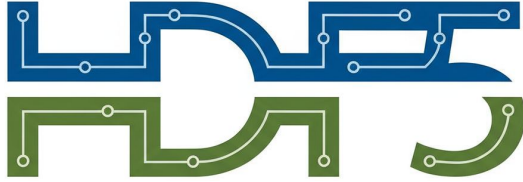metadata

group

dataset
metadata

metadata

# Cloud-ready HDF5

Cloud-ready HDF5 encompasses two complementary approaches for efficiently working with HDF5 data in cloud environments:

- Cloud-Optimized HDF5 (COH5) is a practice that restructures HDF5 files to enable efficient partial reading from S3-compatible object stores.
- HDF5 Scalable Data Service (HSDS) is a REST-based service that provides HDF5 functionality through HTTP APIs.
  - Stores HDF5 data in an object storage-friendly format
  - Multi-reader/multi-writer (MRMW) support
  - Supports horizontal scaling across multiple service instances
  - Offers fine-grained access control and authentication

COH5 focuses on how the data is structured at rest, HSDS provides the mechanism for dynamic access and scaling in a cloud environment. HSDS can leverage COH5 files for even better performance, but it can also work with non-optimized HDF5 files by intelligently managing their access from object storage.

# AI/ML-ready HDF5



**Does going to the doctor make you sick?**

| Name | Visited Doctor | Sick One Week Later |
|------|------|------|
| Alice | Yes | Sick |
| Bob | Yes | Sick |
| Carol | Yes | Sick |
| David | Yes | Healthy |
| Emily | No | Healthy |
| Frank | No | Healthy |
| Grace | No | Healthy |
| Henry | No | Sick |

*AI-ready HDF5* refers to HDF5 containers that are structured using HDF5's hierarchical model and strategically designed, annotated, and managed to meet the specific requirements of AI/ML workflows. This includes optimizing for efficient data access and processing by AI algorithms, ensuring data quality and integrity, providing comprehensive and standardized metadata for <u>context</u> and reproducibility, facilitating seamless integration with AI tools and frameworks, and adhering to FAIR data principles to maximize utility and collaboration.

# Community-ready HDF5

An evolution of the HDF5 ecosystem that is structured, governed, and developed in a way that actively invites, enables, and values contributions from a broad and diverse community of users, developers, and stakeholders. It goes beyond open source to embody principles of openness, collaboration, transparency, and sustainability. A community-ready HDF5 would demonstrate the following characteristics:

1. Transparent and Inclusive Governance
2. Contribution-Friendly Infrastructure
3. Open and Responsive Development Process
4. Roadmap Shaped by Community Needs
5. Support for Ecosystem Collaboration
6. Accessible Education and Outreach
7. Trustworthy and Accountable

**This is THE challenge!**



https://github.com/HDFAlliance

# Parting Thoughts

Language is a tool that builds worlds - or breaks them.

HDF5's ethos is exploratory and forward-looking, one of humility, openness, and readiness for adaptation.

HDF5 is more than a passive data container but rather an active mediator of scientific communication and knowledge preservation.



Thomas Jefferson Building, interior, Library of Congress, Washington DC