# ARK persistent identifiers come of age: 21 years, 1035 institutions

*John Kunze, Dave Vieglais, Tim Clark*

December 2022

ARK Alliance

arks.org

# Three ARK (Archival Resource Key) Use Cases

1. Physical Samples for Earth Sciences

2. ARKs for Biomedical AI

3. Metadata Terms Crowdsourced



Dave Vieglais
University of Kansas
**vieglais@ku.edu**

Tim Clark
University of Virginia
**twc8q@virginia.edu**

John Kunze
ARK Alliance
**jakkbl@gmail.com**

arks.org

# Why care about ARK identifiers?

Because robust web links are rare – the average URL lifetime is 100 days

ARKs serve as persistent identifiers (PIDs) with metadata
- found in the Data Citation Index, HathiTrust, Wikipedia, Wikidata, Internet Archive, ORCID profiles, etc.

In contrast to other PIDs, ARKs have
- no fees, no limits, no walled gardens (decentralized)
- very flexible metadata, including none
- can be assigned to anything digital, physical, or conceptual

# ARK anatomy

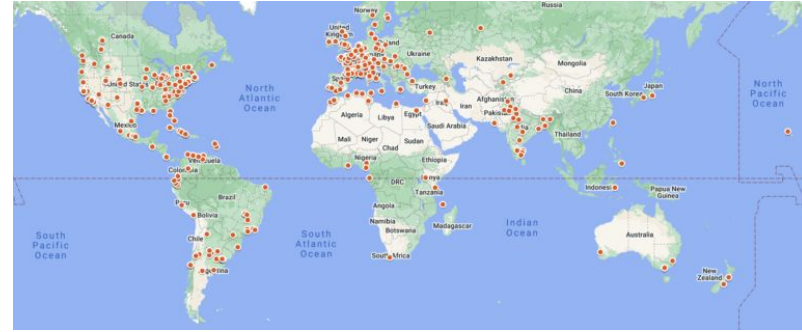A labelled URL with a globally unique identity inside it



https://n2t.net/ark:/12345/fk1234

makes ARK actionable (the resolver)

core globally unique identity (independent of web and hostname)
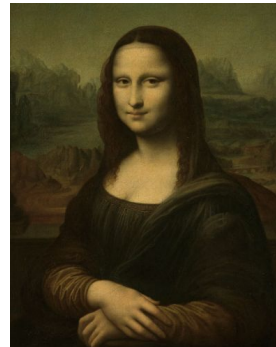
arks.org

# ARK organizations

3.2 billion ARKs created by 1035 institutions –
libraries, archives, museums, publishers,
educators, etc.  For example,



Internet Archive
Bodleian Libraries
Berkeley Law Library
Bibliothèque Mazarine
New York Public Library
French National Archives
National Library of Austria
Library and Archives Canada

University of California Berkeley
Smithsonian National Museum
National Library of France
University of Chicago
Musée du Louvre
Family Search
British Library
Google

https://n2t.net/ark:/53355/cl010066723 →

# Next up: Dave Vieglais

1. Physical Samples
   for Earth Sciences

Dave Vieglais
University of Kansas
**vieglais@ku.edu**

arks.org

# iSamples: The Internet of Samples

The Internet of Samples (iSamples) is a standards-based collaboration to uniquely, consistently, and conveniently **identify material samples**, **record core metadata** about them, and **link** them to **other samples, data, and research products**.

# Factual Basis

- Physical samples typically gathered for analysis
- Analysis may be immediate or many years later
- Samples may be consumed in analysis or preserved in perpetuity
- Derived products such as analyses and publication form graphs of knowledge
- Often beneficial to verify analysis, traversing from publication to original samples
- Linking by identifiers is also essential for attribution, ideally transitive from publications to original sample collector
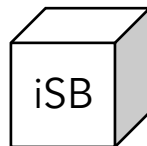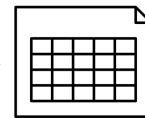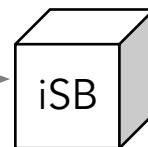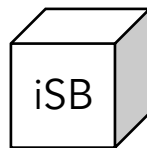
# iSamples Components



- Adapters transform to core records
- Adapters, models for vocabulary mapping
- Records advertised using sitemaps
- Records exposed as JSON
- Records indexed for discovery
- Identifiers are key element

Sample → metadata → Collection → transform → iSB → iSB → 

Model

"Central"
Collates multiple sources

Linked Resources

iSamples

arks.org

# iSamples Summary

- Globally unique, resolvable identifiers are a critical core attribute of physical samples.

- These identifiers facilitate the rapid development of power linked data environments.

- ARK identifiers fill this role with a combination of technical characteristics, reliability, and accessibility.



iSamples

arks.org

# Next up: Tim Clark

2. ARKs for
Biomedical AI

Tim Clark
University of Virginia
**twc8q@virginia.edu**

arks.org
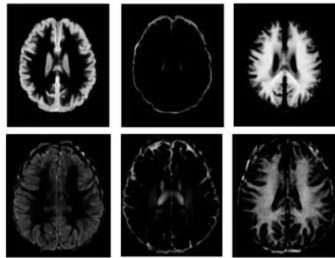
AI is transforming biomedicine...

arks.org

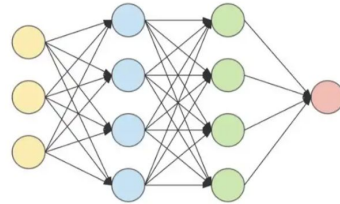...by making very complex predictions...

arks.org

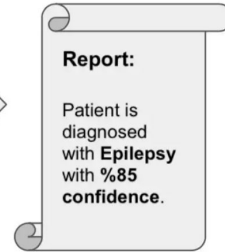...which require explanation and robust evidence!

arks.org

**Epilepsy Detection Model with Brain MRI Data**
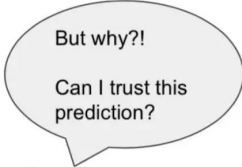
Brain MRI data

Complex ML model

Report:

Patient is diagnosed with **Epilepsy** with **%85 confidence**.

But why?!

Can I trust this prediction?

problem we faced in the previous example. Here are some of the questions

arks.org

*Why did you make that prediction?*

*Can I trust the AI models?*

# Two Biomedical Use Cases
# for Deep AI/ML

1. Predictive model of cellular response to drugs based on very high-dimensional research data and deep learning.

2. Predictive model notifying physicians 7 days ahead of potential life-threatening events on infants in NICU.

arks.org

# Results Validation and Explainability

- Evidence graphs describe how results were obtained & provide supporting evidence for correctness of results.

- We give graph node an ARK resolvable to supporting data, software, and computational parameters- explainability!

- ARKs have huge advantage where thousands of chained computations are performed, and provide open metadata.

arks.org

# 1. Predictive response modeling of normal & diseased human cells

- Construct an accurate cellular component architecture based on very high-dimensional research data and AI.

- Accurately predict cellular response to biochemical perturbation (e.g. new drugs, etc.) using deep learning.

- Be able to interpret and explain model results robustly.

arks.org

adapted from Qin et al. 2021 - A multi-scale map of cell structure fusing protein images and interactions. *Nature* 600:536-542. 16 December 2021. https://doi.org/10.1038/s41586-021-04115-9

# CM4AI FAIR Integration Model



**FAIRSCAPE integration with CM4AI Tools and Data pipelines** (adapted from CM4AI Data Dictionary Requirements Version 1.0). FAIRSCAPE holds comprehensive FAIR metadata and deep provenance graphs on all CM4AI data objects. Data is initially held locally, and accessed through contentURIs in the CM4AI metadata.

arks.org

# Preliminary Map of Cell Architecture
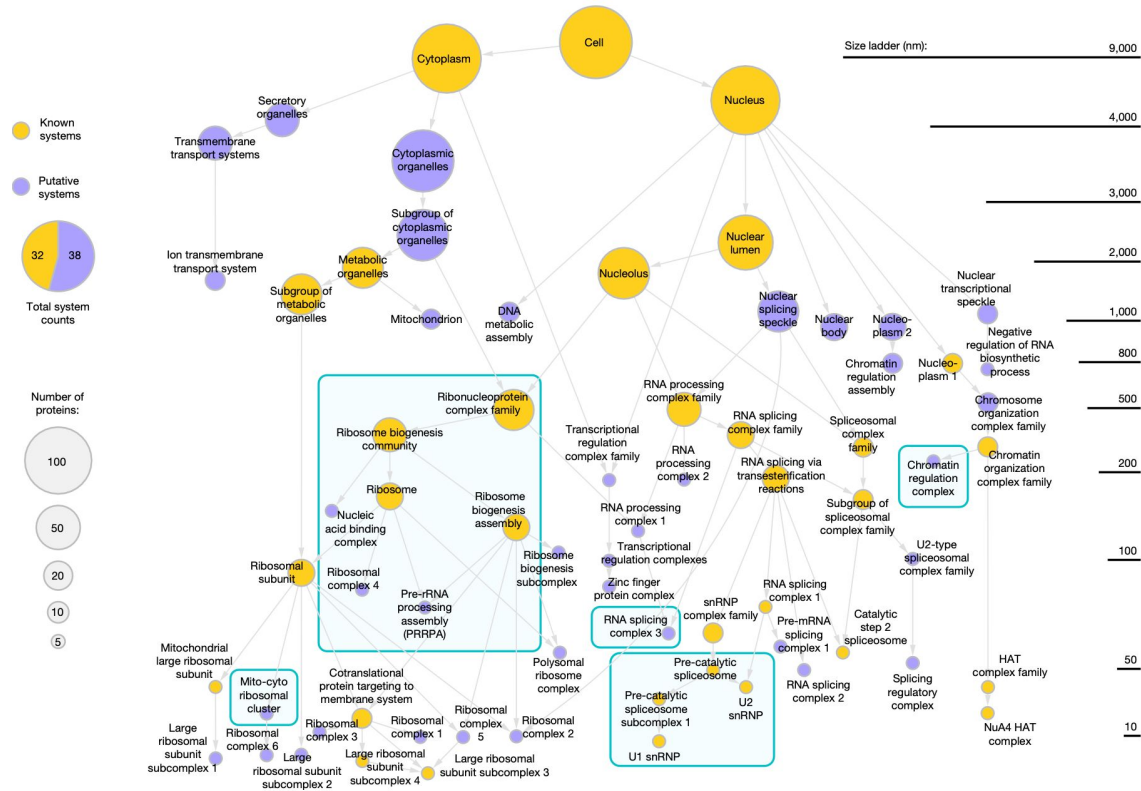


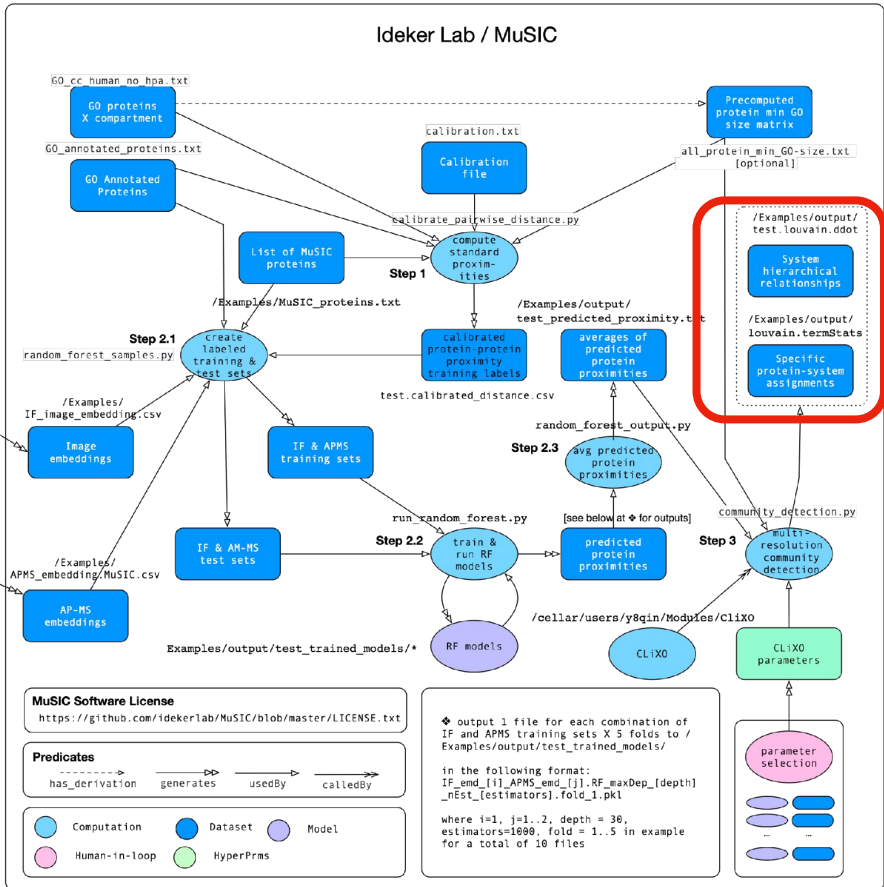adapted from Qin et al. 2021 - A multi-scale map of cell structure fusing protein images and interactions. *Nature* 600:536-542. 16 December 2021. https://doi.org/10.1038/s41586-021-04115-9

adapted from Clark & Al Manor et al. 2023. CM4AI MuSIC 1.0 Provenance Map. NIH Bridge2AI program, CM4AI Deliverable 1.3.

# FAIRSCAPE: initial use case



- NICU time series adverse event prediction on ~ 6,000 babies, > 100 algorithms, ~ 3,000 candidate features, 100 TB data

- Algorithmic result clustering > 17,000 computations per result

[7] J. Niestroy et al., "Discovery of signatures of fatal neonatal illness in vital signs using highly comparative time-series analysis," (2022) *npj Digital Medicine*. preprint: https://doi.org/10.1101/2021.03.26.437138.

arks.org

# 2. Predictive Analytics in the Neonatal ICU

- Predictive analytics: neonatal ICU @ UVA Med Center

- 6,000 NICU babies, 10 years of vital signs data

- 80 time series algorithms X dozens of parameter sets

- Goal: Predict adverse medical events 7 days in advance

arks.org

# Single Patient, 7-node Graphlet for HCTSA Computations



(complete graph for results > 17,000 nodes)

# JSON-LD, <u>schema.org</u> & EVI Evidence Graphs

```
In [14]:  output_id  = job_data['Output Identifiers'][1]
          r = requests.get('http://ors.uvadcos.io/' + output_id)
          output_meta = r.json()

In [15]:  RenderJSON(output_meta)
```

```
              "@context": ⊕{2 items},
              "@id": "ark:99999/1a22ea0a-2b15-4515-a851-6d2b2f7db211",
              "@type": "Dataset",
              "distribution": ⊕[1 item],
              "eg:evidenceGraph": ⊖{
                  "@id": "ark:99999/1a22ea0a-2b15-4515-a851-6d2b2f7db211",
                  "@type": "Dataset",
                  "distribution": ⊕{6 items},
                  "eg:generatedBy": ⊖{
                      "@id": "ark:99999/23eb4cfe-91d1-464b-8e4f-7a372724547f",
                      "@type": "eg:Computation",
                      "dateEnded": "Wednesday, October 30, 2019 07:16:15",
                      "dateStarted": "Wednesday, October 30, 2019 07:15:58",
                      "eg:usedDataset": ⊖{
                          "@id": "ark:99999/9f627f3e-c59f-495f-b2f5-13d1a06a622a",
                          "author": ⊕{3 items},
                          "description": "Heart Rate Measures from patient from admission to discharge.",
                          "eg:generatedBy": ⊕{5 items},
                          "name": "Patient 7129 HR"
                      },
                      "eg:usedSoftware": ⊕{6 items}
                  },
                  "name": "Output from job Job 7294"
              },
              "eg:generatedBy": ⊖{
                  "@id": "ark:99999/23eb4cfe-91d1-464b-8e4f-7a372724547f"
              },
              "identifierStatus": "DRAFT",
              "name": "Output from job Job 7294",
              "sdPublicationDate": "2019-10-30T19:16:29.76Z",
```

**return metadata - formal ontology terms**

**DAG of provenance / evidence**

**derived by this computation**

**from this dataset**

29

# Why do we use ARKs?

- Persistent IDs for data (the evidence graph nodes)

- Free to mint ARKs and their metadata is flexible

- Large ecosystem of users & developers

arks.org

# We used ARKs for

- Complex evidence graph on AI/ML predictions

- Every node (dataset, computation, software) resolvable

- Each node persistently identified with an ARK

arks.org

# Conclusion

- ARKs are a useful, flexible, scalable persistent ID model.

- Especially useful for traceable complex computations.

- We are happy to chat with prospective ARK users.

arks.org

# Next up: John Kunze

3. Metadata Terms
Crowdsourced

John Kunze
ARK Alliance
**jakkbl@gmail.com**

arks.org

# Yamz.net and ARKs for metadata terms

Vocabulary creation, sharing, and standards – better, faster, cheaper

*John Kunze*

<MRC>

Drexel
UNIVERSITY
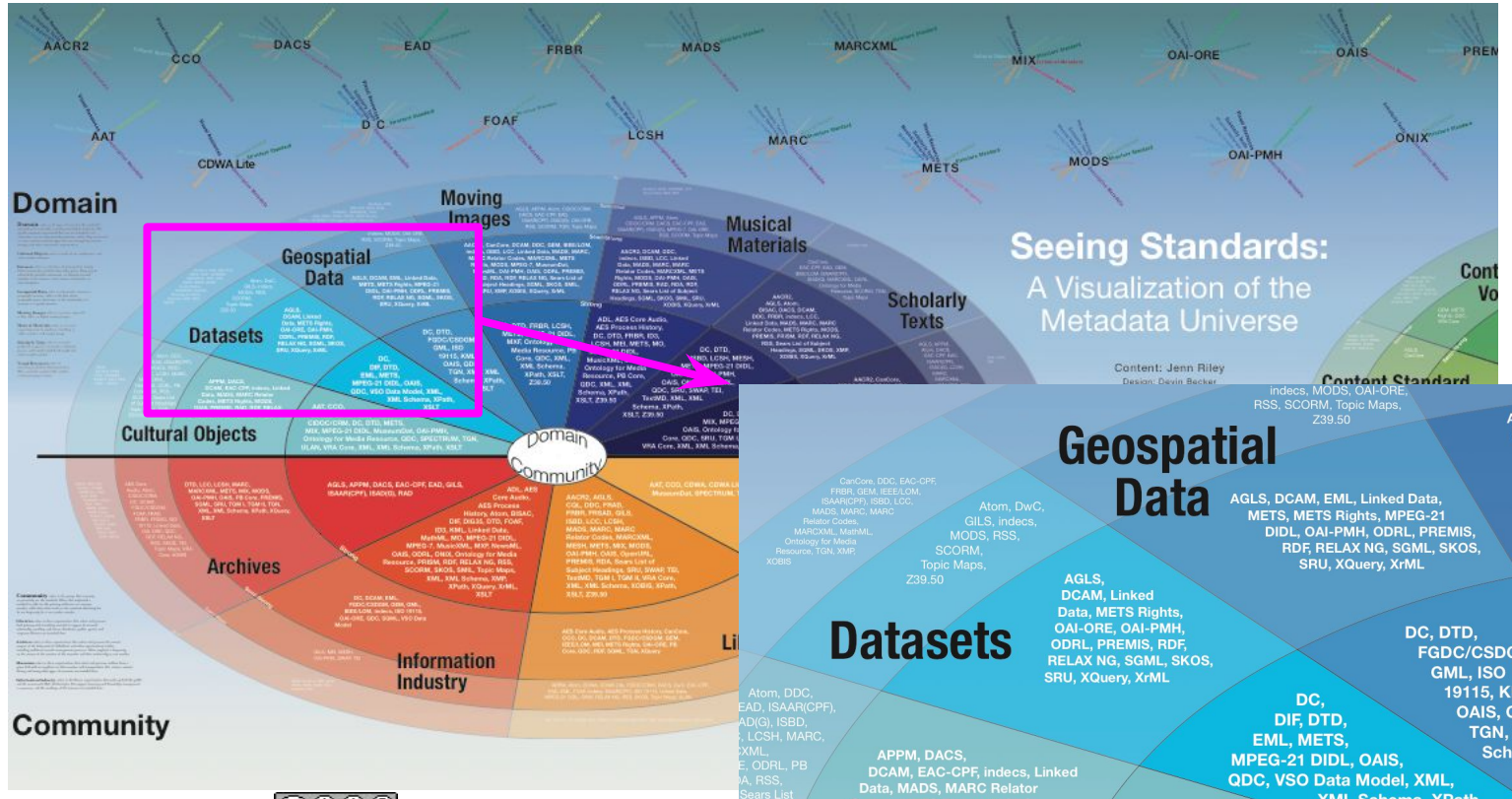
arks.org

# Standardized metadata is rare



**Theory**

- Our boss: "*We use Dublin Core, PREMIS, X, Y, and Z, for interoperability.*"

**Practice**

- Cataloger/Archivist/Scientist: "*Frankly, we use those vocabularies, but with our own local modifications to make them work for our objects.*"
- Many unofficial dialects – per institution, per laboratory, per project
- Poor interoperability (see *Metadata's Bitter Harvest*, Library Journal, 2004)

arks.org

# The Metadata Universe



**Seeing Standards:**
A Visualization of the Metadata Universe

Content: Jenn Riley
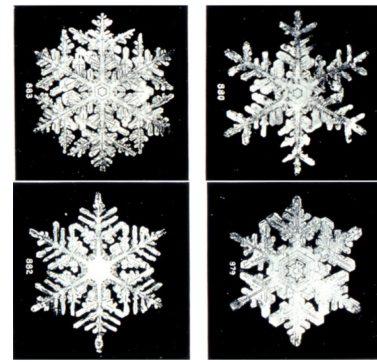Design: Devin Becker

Jenn Riley, IU

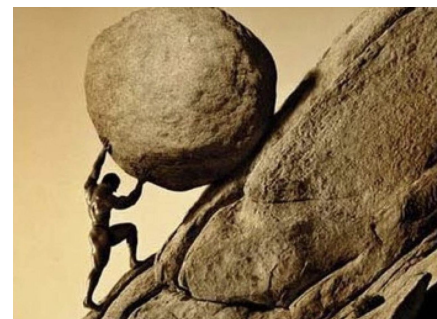arks.org

# Domain dialects – similar but different

**Example**: Earth Science > Cryospheric (frozen water) Science

- 28 different definitions of "glacier"
- 8 different definitions of "puddle"
- 13 different definitions of "firn" (old snow)
- 10 different definitions of "frazil ice" (fine spicules of floating ice)
- 7 different definitions of "ogive" (bands of light and dark ice in a glacier)
- … and on and on

Sound familiar? What about your domain?

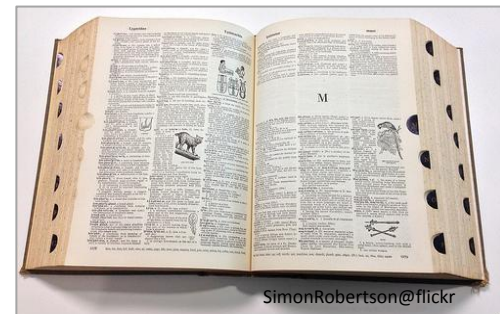# Which term definitions to keep or toss?



Traditional approach
- Busy experts in *<your fast-moving field>* on metadata standards committees
- … often take years to reach consensus
- … often no field testing is logistically feasible

An alternate approach is Yamz.net.

arks.org

# Yamz.net (Yet another metadata (zoo))


SimonRobertson@flickr

Yamz is not a standard, nor an ontology

- Yamz is a *living* dictionary of metadata terms
- Each term gets an ARK permalink, becoming a kind of
- … proposed nano-standard – some are upvoted, others not
- Reputation-based voting (like Stack Overflow) helps choosing

Yamz is a *microservice* for sharing, testing, and revising metadata dialect

- All parts of metadata "speech", all domains, all ontologies
- Field testing via practitioners and reputation-based voting
- Ontologies, software, and predicates reference metadata terms via ARKs

arks.org

# Please reach out with feedback

1. Physical Samples
for Earth Sciences

2. ARKs for
Biomedical AI

3. Metadata Terms
Crowdsourced

Dave Vieglais
University of Kansas
**vieglais@ku.edu**

Tim Clark
University of Virginia
**twc8q@virginia.edu**

John Kunze
ARK Alliance
**jakkbl@gmail.com**

arks.org