

Capturing provenance throughout the biodiversity genomics pipeline with RO- Crate

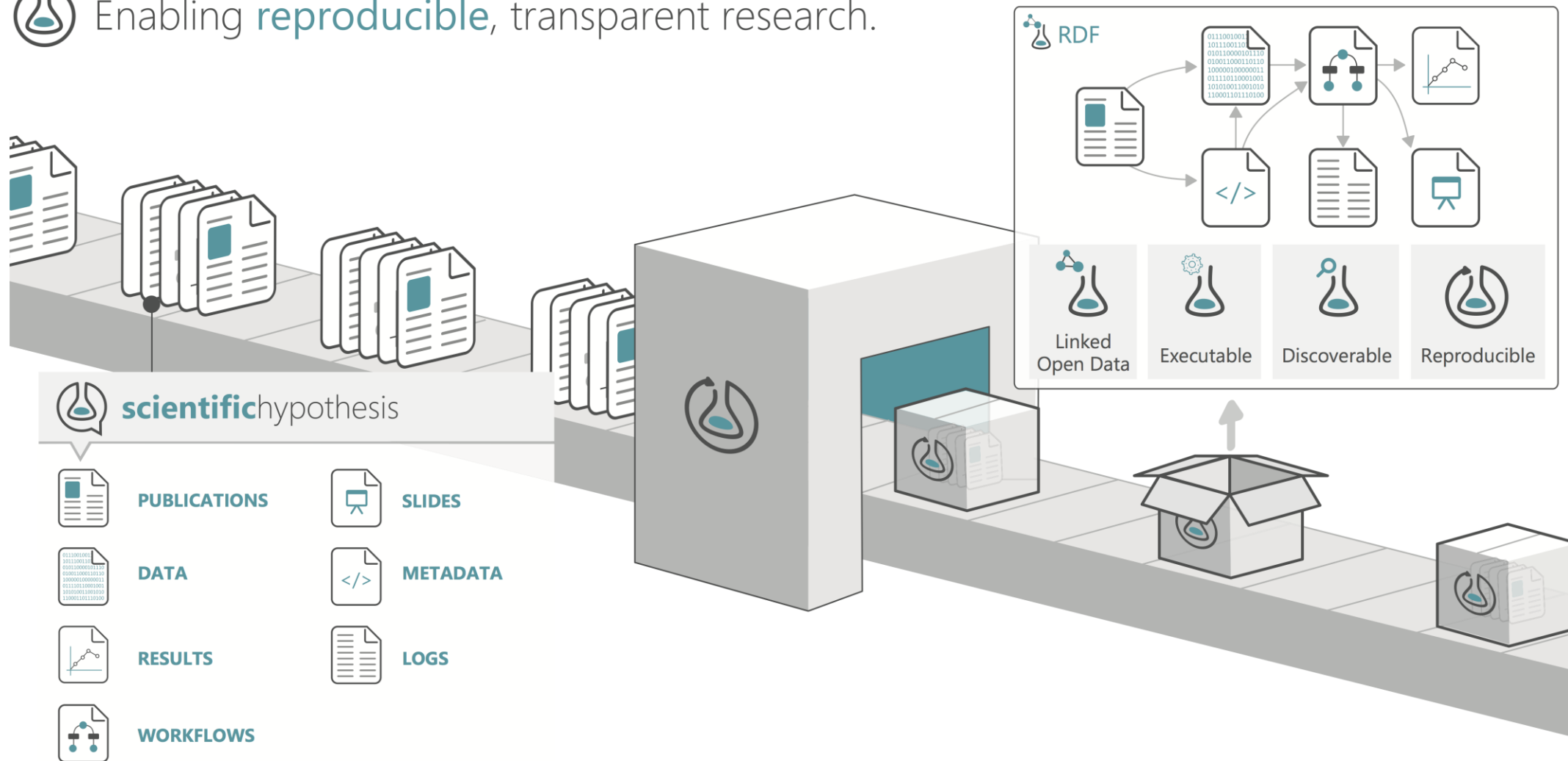
Eli Chadwick

University of Manchester



What is RO-Crate?

 Enabling **reproducible**, transparent research.



What is RO-Crate?

 Enabling **reproducible**, transparent research.

The “who, what, where, how” of your research



RESULTS



LOGS



WORKFLOWS



RDF



Benefits of RO-Crate

Easy to include

Add metadata file next to existing data; distribute them together

Easy to consume

Standardised, interoperable format

Works with your data

A general base, with custom domain profiles

Uses *linked data*

Naturally interconnected with other RO-Crates and the wider web

RO-Crate Profiles

Extensions to RO-Crate to support needs of a particular community.
For example:



Description of a workflow
execution

- inputs, outputs, who ran it, when, environment...



Describing life science experiments

- ISA (Investigation, Study, Assay)
- lab protocols & processes



Why use RO-Crate for biodiversity genomics?

- ERGA: PDF genome assembly reports list associated records from ENA, BioSamples, ...
- RO-Crate object links to associated records directly – accessions, quality metrics, workflows
- More explorable, more machine actionable

ERGA Assembly Report
v24.04.03_beta

Tags: ERGA-BGE

TxID	1464561
ToLID	idCulLati1
Species	Culex laticinctus
Class	Insecta
Order	Diptera

Genome Traits	Expected	Observed
Haploid size (bp)	726,182,214	833,812,495
Haploid Number	3 (source: ancestor)	3
Ploidy	3 (source: ancestor)	2
Sample Sex	XX	XX

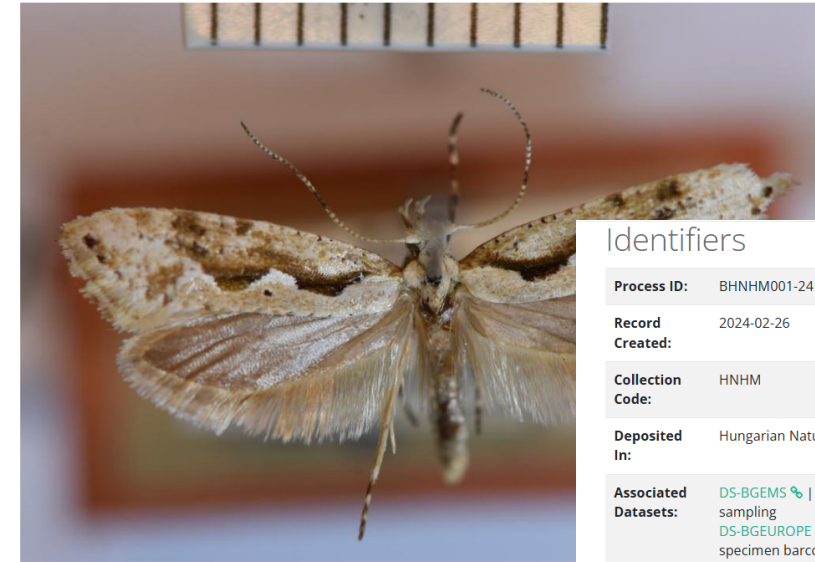
EBP metrics summary and curation notes

Obtained EBP quality metric for prj: E-BGE

Why use RO-Crate for biodiversity genomics?

- BOLD: does things completely differently – single record covers sample and processes
- Less metadata about sub-processes – sequencing vs genome assembly
- Can we move toward a more unified approach?
-> RO-Crate

Specimen Images



Identifiers

Process ID:	BHNM001-24	Sample ID:	BGE_00146_A01
Record Created:	2024-02-26	Museum ID:	HNHM-LEP-11332
Collection Code:	HNHM	Field ID:	
Deposited In:	Hungarian Natural History Museum		
Associated Datasets:	DS-BGEMS BGE Biodiversity Genomics Europe: Museum sampling DS-BGEUROPE BGE Biodiversity Genomics Europe: Curated specimen barcoding output		

Specimen Linkout:

Checksum: ae8a513c896e36dc0c17025a8479972c



Taxonomy

Applying RO-Crate to BGE

Computational
analysis

Real-world
processes (wet lab,
sample collection)

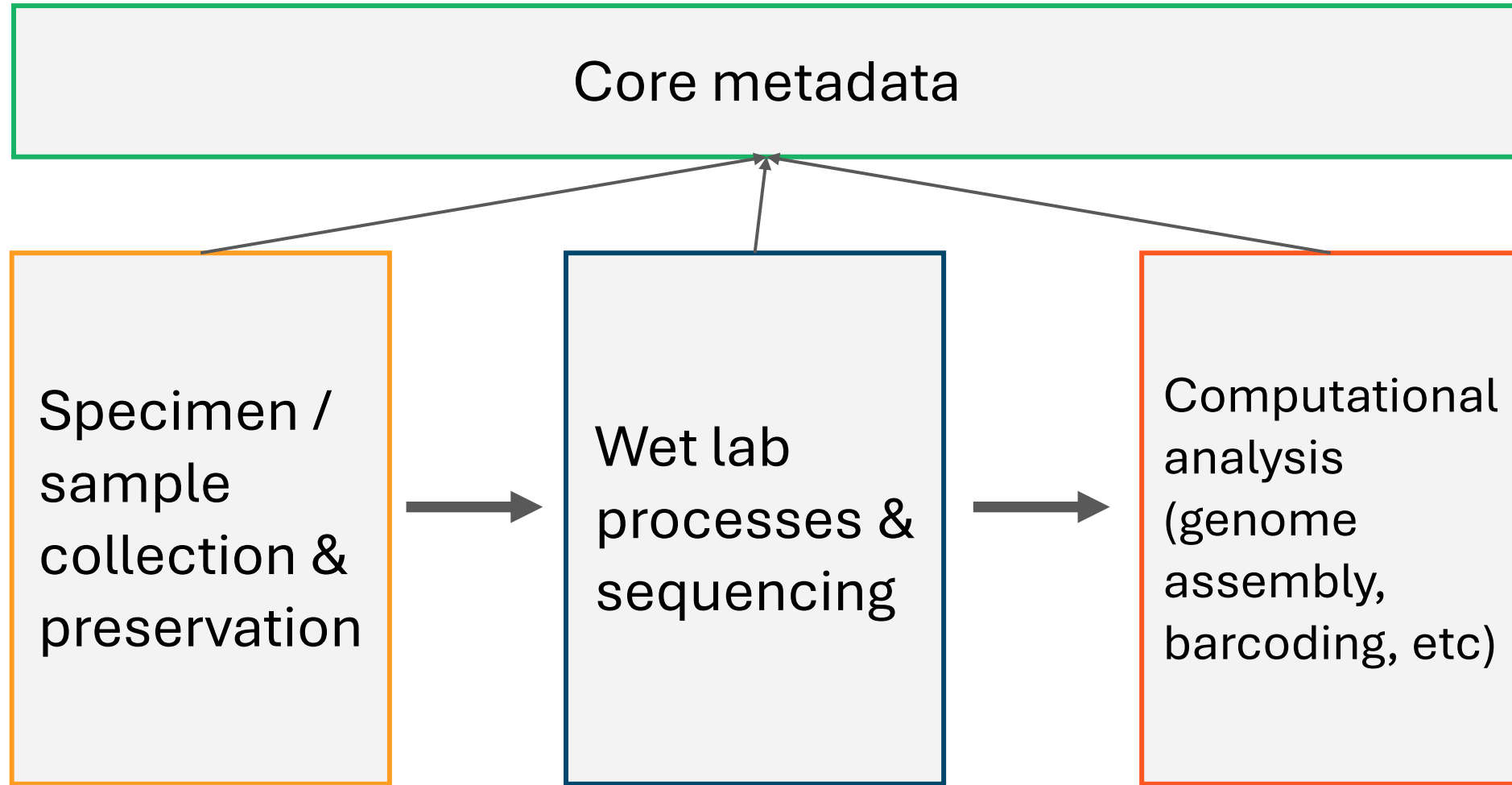
*Goal: full
provenance from
sample to barcode*

Genomics-
specific metadata

ERGA-specific
metadata

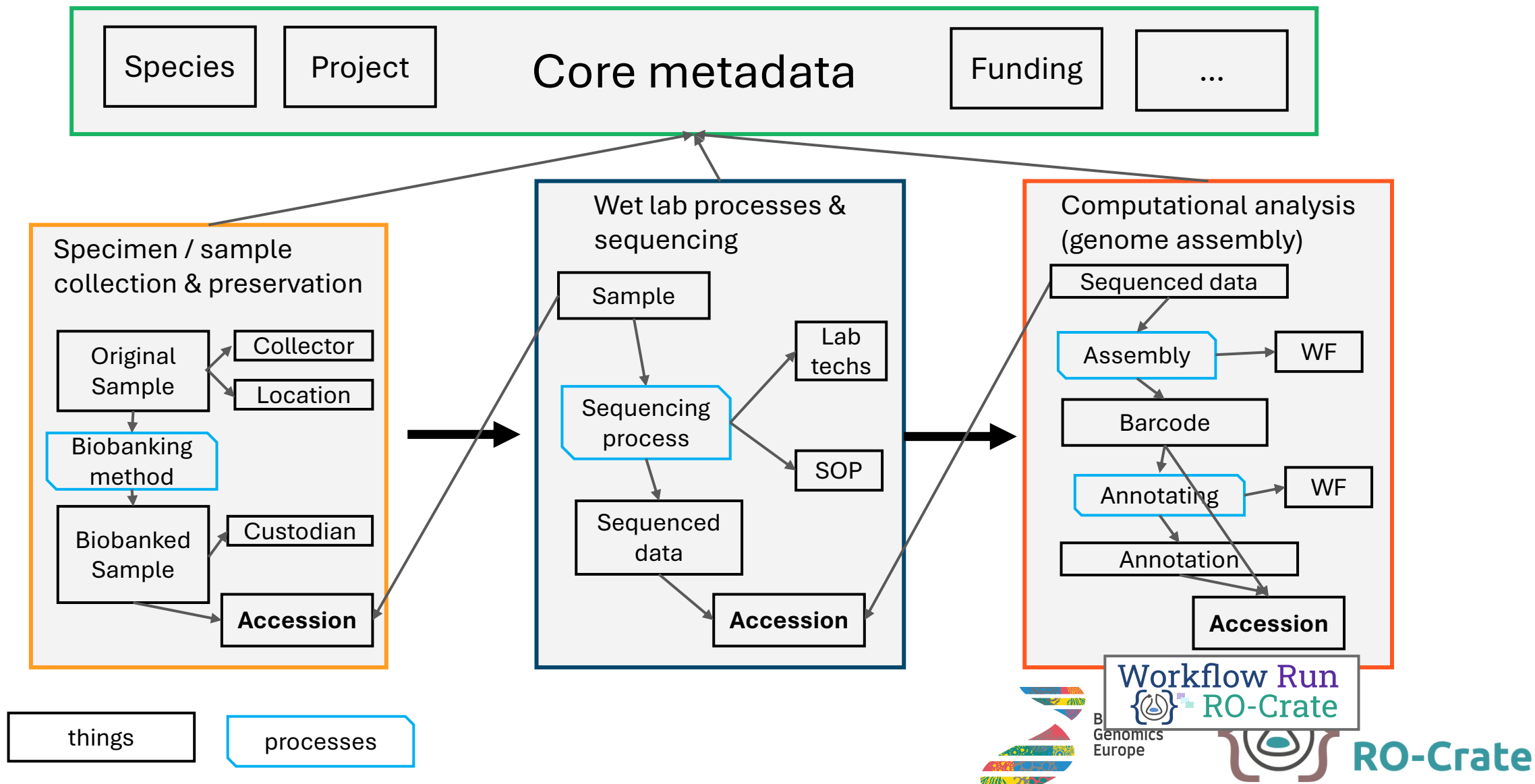
iBOL-specific
metadata

BGE RO-Crate Profile



A sequence of processes

BGE RO-Crate Profile



BOLD example

Taxonomy

Kingdom:	Animalia	Subfamily:	Hesperiinae
Phylum:	Arthropoda	Tribe:	Calpodini

Core metadata

Specimen /
sample
collection &
preservation

Process ID: MHMXN361-07

Sample ID: 07-SRNP-32653

lab
processes &
sequencing

Computational
analysis
(genome

Sequence: COI-5P

Sequence ID: MHMXN361-07.COI-5P

GenBank
Accession: JF761761

Primers
Forward: LepF1 (ATTCAACCAATCATAAGATATTGG)
MLepF1 (GCTTTCCACGAATAAATAATA)

Primers
Reverse: LepR1 (TAAACTTCTGGATGTCCAAAAATCA)
MLepR1 (CCTGTCCAGCTCCATTTT)

Sequence Run
Site: Centre for Biodiversity Genomics

TAGGAACATCATTAAAGATTATTAATCCGAACAGAATTAGGAAACCCAGGATCTTTAATTGGAGATGATCAAATTTACAATACTATCGTTACTGCTCATGCTTTTA
TTATAATTTTTTTATAGTAATACCTATTATAATTGGAGGATTTGGAAATGATTAAATCCCTTAATATTAGGGGCTCCCGACATAGCTTTCCCGAATAAACA
ACATAAGATTTTGAATATTACCCCATCTTTAACTCTTTAATTTCAAGAAGAATTGTAGAAAAATGGTCCGGAACAGGTTGAACGTGTTACCCCCCTTTTCAT
CTAATATTGCCCATCAAGGATCTTCGGTCGATTAGCAATTTTTCTTACATTTAGCTGGTATTTCTTCAATCTTAGGGGCTATCAATTTTATTACTACAATTA
TTAATATACGAATAAAACTTATCATTTGATCAAATACCTTTATTTATTTGATCTGTAGGAATTACAGCACTATTATTACTCTTATCTTTACCGTATTAGCTG
CGCTATTACTATACTTTAACTGATCGAAATTTAAATACCTCTTTTTCGACCTGCGGGAGGGGG

A sequence of pro

Collection

Country/Ocean:	Costa Rica (CR)	Collection Date Start:	2007-08-02
----------------	-----------------	---------------------------	------------

Province/State: Guanacaste Province

Specimen

Voucher Status:	Reproduction:	S
-----------------	---------------	---

Tissue Descriptor: Sex:

Specimen Images



Challenges

- Barcodes from BGE don't have GenBank accessions yet
 - BOLD identifiers are less universal – less linkable
 - ENA integration in the works
- “Missing” metadata about sequencing/assembly processes
 - Sequencing protocol(s)
 - Barcode assembly workflow(s)
 - Who did the assembly?
 - Does this metadata live somewhere but it's not in BOLD?
- Creating RO-Crates *during* these processes – no identifiers

Next Steps

- Combined barcode/genome use cases?
- Alignment with other standards ([ISA](#), [Common Provenance Model](#), [Annotated Research Context](#), [GBIF](#), ...)
- Describing other processes:
 - Barcode validation
 - Sample collection & identification against BOLD
- Can YOU pilot this profile and make your data more FAIR?

Thank you!

Learn more: <https://esciencelab.org.uk/bge-ro-crate-profile/>

Get in touch: eli.chadwick@manchester.ac.uk

Come to an RO-Crate drop-in: <https://s.apache.org/ro-crate-regional>

next one: this Wednesday at 15:00 UTC

Thanks for discussions: Joana Pauperio, Stian Soiland-Reyes, Tom Brown, Peter Woolland, Rutger Vos, Nick Juty, Pete Hollingsworth

