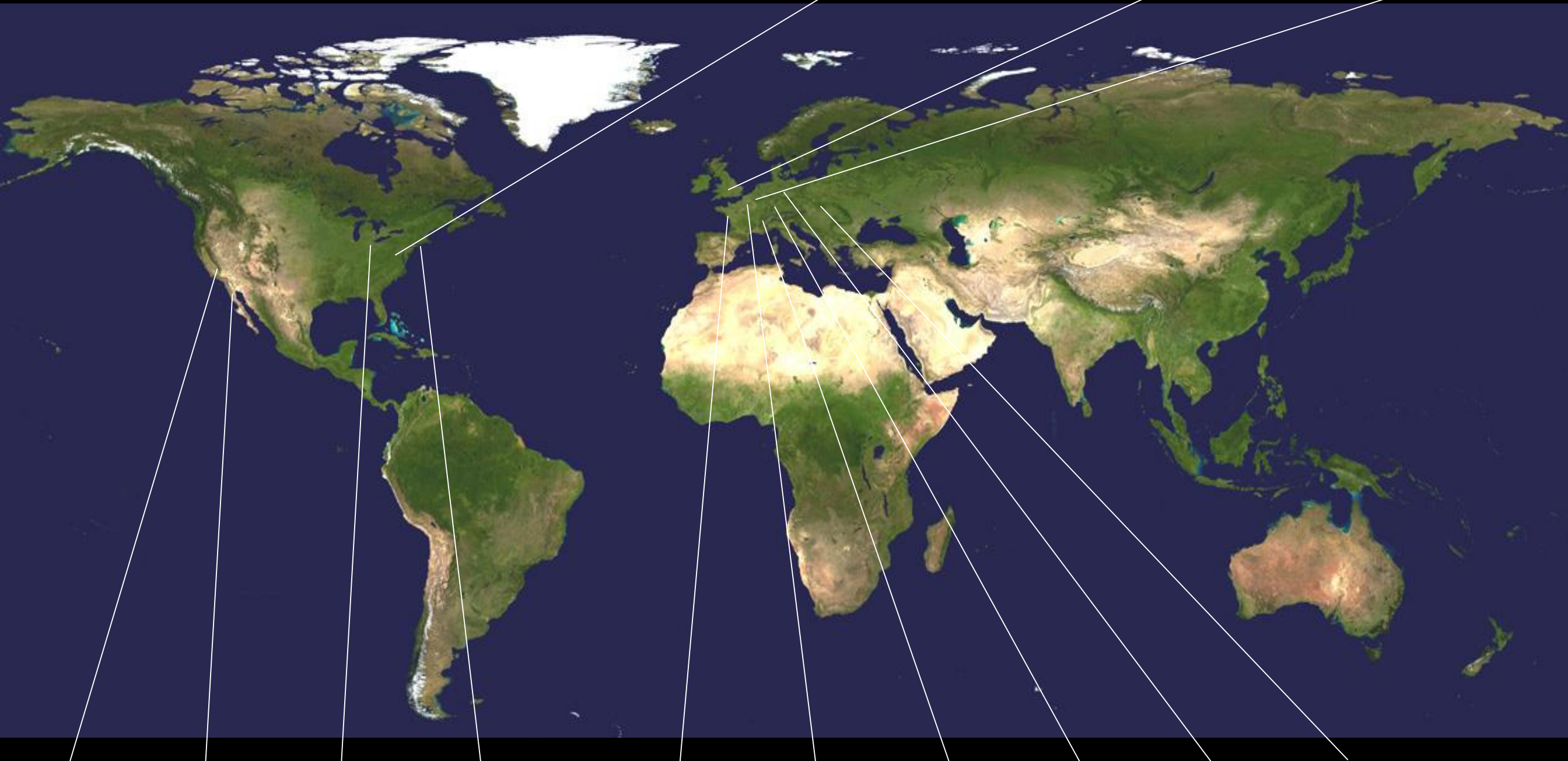# No specimen left behind:
## The NHM Digital Collection Programme

Vince Smith, Natural History Museum, London
Digitisation of Biological Collections, 13 April 2014

1.5-3 BILLION SPECIMENS
1.9 million species
300 years of collection

Washington 125M

London 80M

Paris 60M

San Francisco 28M

Los Angeles 35M

Chicago 25M

New York 30M

Brussels 37M

Leiden 37M

Vienna 35M

Frankfurt 40M
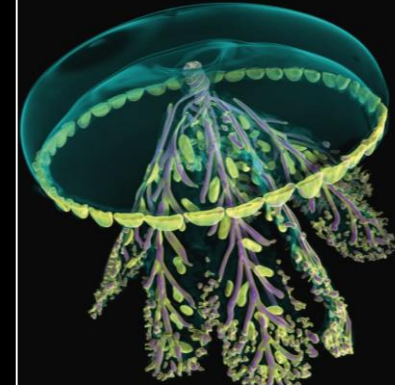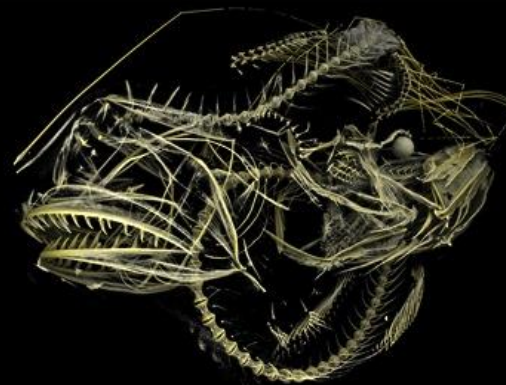
Berlin 30M

St P'berg 32M

Digital surrogates of specimens:
- 3D models and printing
- CT imaging of internal structures
- Atomic, physical and chemical analysis

Using digitisation to enhance public engagement with natural history specimens:
- Augmented reality in gallery
- Animated content for pre- and post-visit apps

An inordinate fondness for beetles:
400,000 species described so far, 100,000 types among the NHM's 10M beetle specimens

# NHM collection

| Collection area | No of objects | No of type specimens | Physical register | Digital data |
|---|---|---|---|---|
| Palaeontology | 6,919,207 | 43,146 | 2,364,232 | 340,636 |
| Mineralogy | 423,563 | 615 | 425,000 | 402,727 |
| Botany | 5,863,000 | 172,750 | 127,200 | 645,222 |
| Entomology | 33,753,257 | 612,796 | 57,197 | 255,000 |
| Zoology | 27,501,350 | 325,000 | 1,986,000 | 1,160,216 |
| Library & archives | 5,460,000 | - | - | - |
| TOTAL | 79,920,377 | 1,154,307 | 4,959,629 | 2,803,801 |



**<3% of NHM specimens are digitised, & even fewer are 'computable'**

# Digital NHM: a strategic priority

**Museum Strategy 2015-2020**

**Thee big narratives**

① Origins and evolution

② Living diversity

③ Sustainable futures

**Four priority areas**

① Digital

② National

③ International

④ London

**Priority Area #1: Digital**

| Objectives | Key projects |
|---|---|
| **Big, open data** | Mass collection digitisation*<br>NHM Data portal*<br>Crowd-sourcing of data* |
| **Global communities** | New website*<br>Social media channels<br>Membership |
| **Innovative platforms** | Location-aware Wi-Fi*<br>Museum App*<br>In-gallery technology |

# iCollections: 2013-2015

- A pilot for the Digital Collections Programme
- iCollections digitisation criteria:
    - Entire collection
    - No existing digitisation pipeline
        *(pinned, slide & herbarium specimens)*
    - High research potential
        *(phenology, morphometrics, migration patterns,
        pest associations, automated species recognition)*
    - Curation opportunity


- Data outputs
    - Image(s)
    - Metadata (what, when & where)
    - Georeference point localities

**NHM iCollections project:**
- UK butterflies & moths
- 500k specimens
- 2 mins per specimen
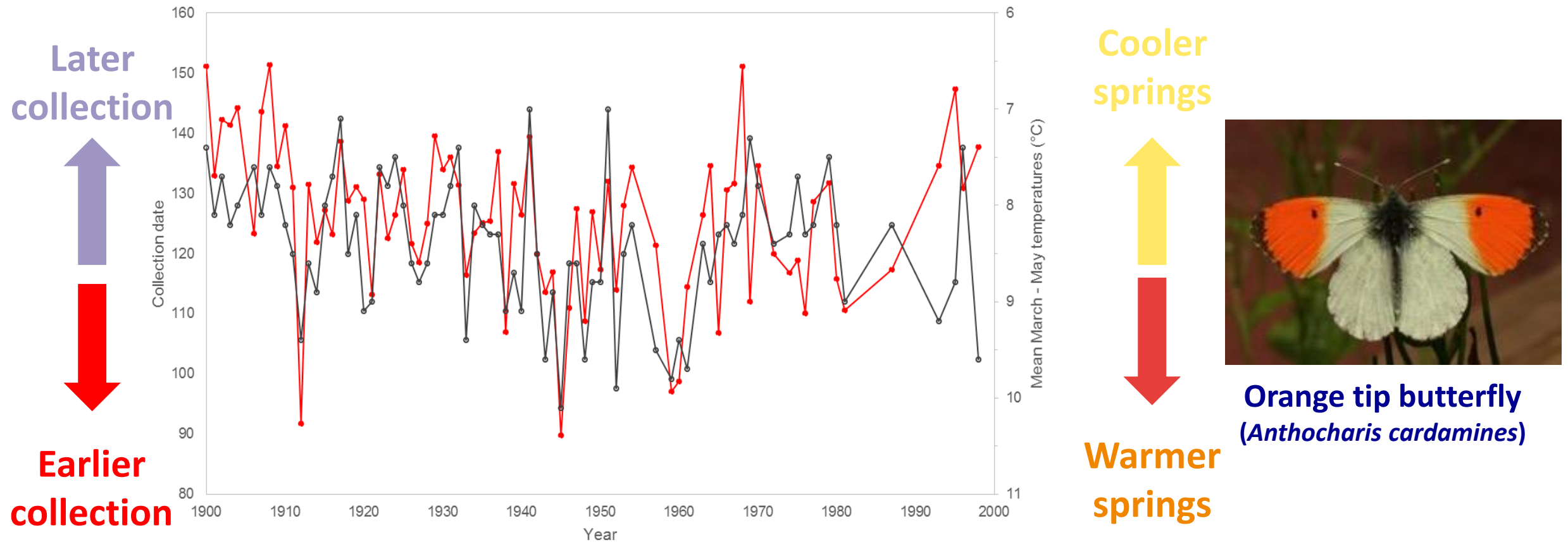- £1 per specimen

**Large-scale digitisation:**
- High-throughput digitisation workflows
- Informatic pipelines
- Computer-assisted object recognition

# iCollections digitisation process

# iCollections research: long-term trends in phenology



**Orange tip butterfly**
(*Anthocharis cardamines*)

- **Many species emerging earlier and earlier each year**
- **Initial collection date & temperature highly correlated**
- **A unique marker on phenological response before recent climate change**
- **Longer time perspective than most observational records (BMS post-1976)**
- **Museum data available for rare or hard to record species**

# NHM Digital Collections Programme

*"To collate, organise and make available to global scientific and public audiences one of the world's most important natural history collections, delivering:*

*- an online specimen- / lot-level data base for all holding*

*- core meta-data and / or images for key parts of the collection, and*

*- flexible informatics and visualisation tools"*

**Target = 20 million specimens digitised in 5 years**

| | 2 year | 5 year | 10 year |
|---|---|---|---|
| **POLICY & PROTOCOL** | Defined data policy and standards | Policies embedded in NHM operating practises | Leaders of process in the digital curatorial world |
| **DATA CAPTURE** | Prioritised digitisation Workflows piloted | Portfolio of mass digitisation output projects | Some major collections digitised |
| **PEOPLE & SKILLS** | Task force formed and operating | Best-practice processes integrated into training | Digital curation as a core part of our practice |
| **INFRASTRUCTURE** | Refined collections database, tools & hardware | Future collections database implemented | Broad connections to other large digital collections |
| **STAKEHOLDERS & GOVERNANCE** | Key user communities engaged | Peer to peer development | Proactive engagement of emerging audiences |
| **PARTNERSHIPS** | Partners involved in pilot projects | Fully funded digitisation portfolio | Major international coalitions |
| **RESEARCH** | Research-orientated projects & initiatives | Collaborative research material published | Major contributions to grand challenges |
| **ACCESS** | Live NHM Data Portal | Tools , visualisations & analytics | Integrated global network of users |

# Digital Collections Task Force

- Core pool of expertise, ideas and capacity for the Museum's digitisation activity

- Key individuals across departments

- 50% or more of their time on projects under DCP

- Individuals seconded to relevant projects

- Coordinated by a steering group made up of units lead individuals and chaired by the Head of Collections Digitisation

- Membership expected to change as the portfolio of projects evolves
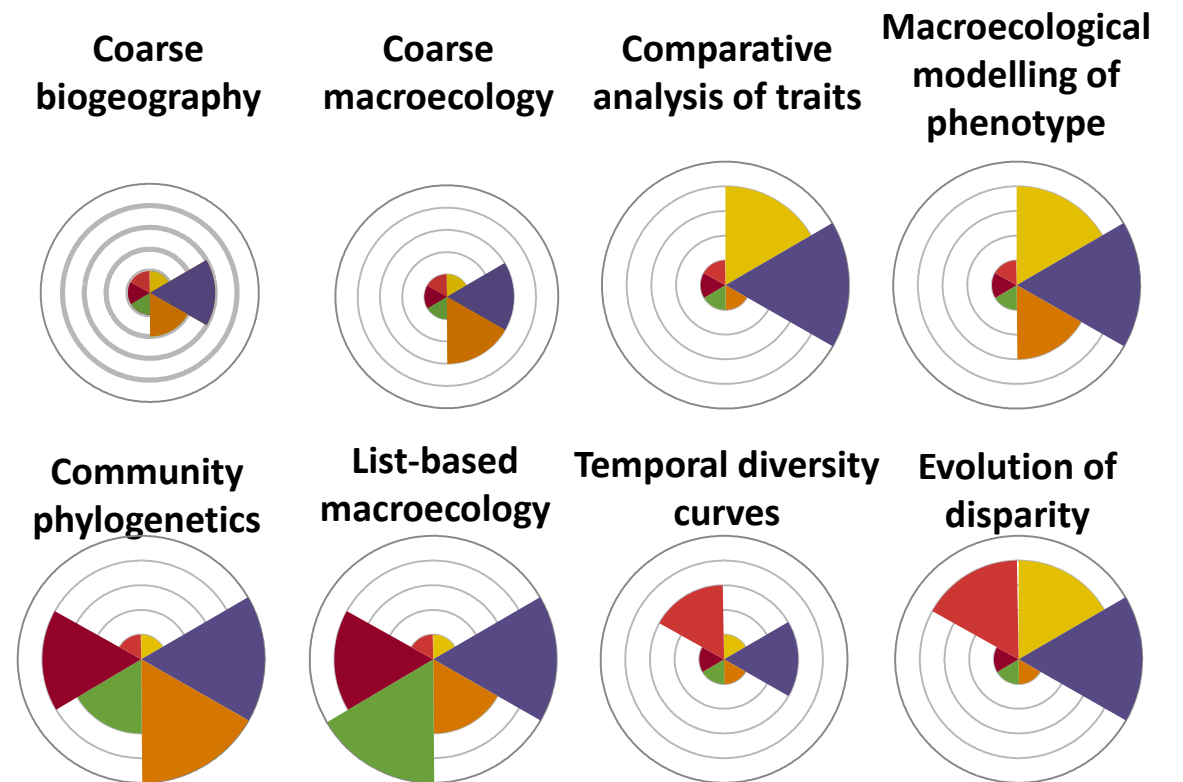
- **Change from within**

# Question-oriented prioritisation

**Core meta-data categories:**



Stratigraphy

Object Image

Collector

Taxonomic name

Date of collection

Geographical location

**What meta-data are required for different types of study?**

Coarse biogeography

Coarse macroecology

Comparative analysis of traits

Macroecological modelling of phenotype

Community phylogenetics

List-based macroecology

Temporal diversity curves

Evolution of disparity

# Digital Collections Programme: pilot projects

## Pilot 1 – Herbarium sheets

- Trial partnership with Kew & Picturae (cf Naturalis)
- Conveyor-based imaging equipment
- Outsourcing of label transcription
- C.70,000 sheets
- January-April 2015



## Pilot 2 – Palaeontology

- British Fossils (Mesozoic vertebrates)
- Standardised photography
- Capturing taxonomic & stratigraphic metadata
- February 2015 – March 2016



## Pilot 3 – Microscopic slides

- High-throughput slide scanner
- Satscan to capture label information
- Call for pilot project ideas
- April 2015 onwards

# NHM & Kew Herbarium sheet digitisation



**33k Specimens per day, 3 shifts (6am-10pm), Netherlands collection complete in 1.5 years**
**€1.29 Euros per specimen image (if outsourced), transcription at similar cost**

# Crowdsourcing

# Crowdsourcing platforms



**Herbarium @Home**
http://herbariaunited.org/atHome/

**Smithsonian Transcription Center**
https://transcription.si.edu/

**Atlas of Living Australia**
http://volunteer.ala.org.au/

**Les Herbonautes**
http://lesherbonautes.mnhn.fr/

**Notes from Nature**
http://www.notesfromnature

**Crowdcrafting**
http://crowdcrafting.org

NATURAL HISTORY MUSEUM

# Notes from Nature (Zooniverse)

## NHM Ornithology Registers (1837 – 1990)

**Progress**

- Total Images: 2,950
- Active Images: 126
- Complete Images: 2,824
- 307,305 transcriptions
- Circa 1/3$^{rd}$ by one person

# Science uncovered 2014: "Crowdsourcing the Collection"

- Demonstrate digitisation process
- Engage the public in transcription
  - Digitise; web publish images; public transcription; data publication
- Dedicated mobile website

# NHM Data Portal

- A platform for deposition and discovery of NHM research & collections data
- Weekly updates from NHM collection database
- Stable, citable (DataCite) identifiers on datasets & records
- Transparent data quality indicators
- Export, import, & web-services available (Linked Open Data coming soon)
- Built using CKAN, with enhanced mapping, statistical and filtering functions

# NHM Open Data Policy

**Open-by-default policy on collections data & research data sets**

- Creative Commons (CC) Zero waiver on data

- CC BY 4.0 licence on images

**Exceptions to default**

- Exceptions require justification, which include:

    - *Commercial value*
    - *Research competitiveness*
    - *Third party rights*
    - *Sensitive collection information*
    - *Donor or funder conditions*
    - *Confidential documents*

**PSI Directive**

- EU legislation on reuse of public sector information, extended to cover museums, libraries and public archives.  To be implemented in July 2015

- General principal that information should be available for re-use for free or at marginal cost, and if reused, made available to all under equal terms.