# GigaByte – A New Workflow for Rapid Dissemination of Datasets & Tools - Bringing Papers to Life

eResearch 2021, Wellington, NZ

(GIGA)bYte
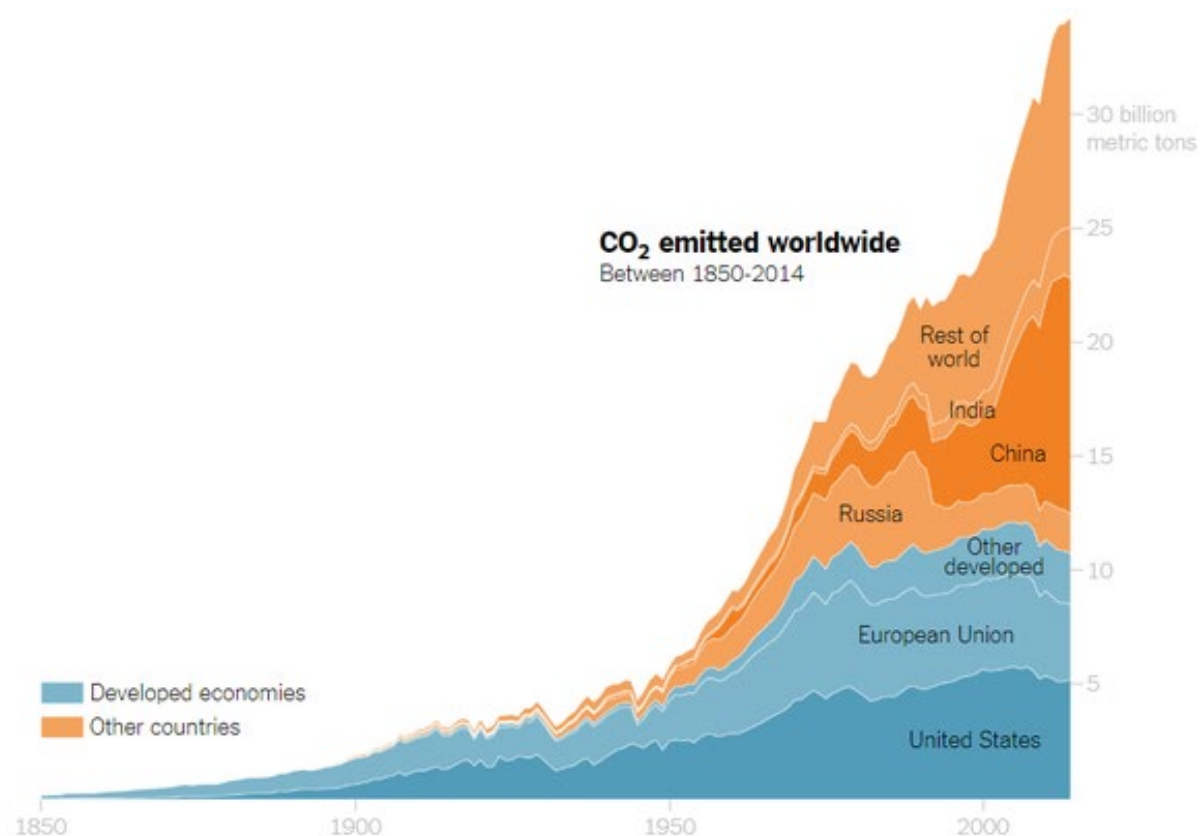Publishing at the Speed of Research

**Nicole Nogoy**

Executive Editor,

ORCID: 0000-0002-5192-9835

(GIGA)ⁿ SCIENCE PRESS

BGI 华大

# Urgent Challenges Research Needs to Address

**Climate change**

**Disease pandemics (e.g., COVID19)**

# Urgent Challenges: Publishing to the Rescue?



- Need to disseminate & communicate information openly to global community
- Expensive. Information is held back by paywalls, APCs, and other barriers

- Need to share this quickly and in a trusted form (peer reviewed)
- Laborious, archaic tech, and untransparent processes

- Research data, software and underlying methods and results need to be shared for scrutiny and re-use
- Hard work and little incentive to share

- Needs to be understandable by policy makers, public
- Barriers of language, jargon, and lack of interaction

# Attempt to Address This: *GigaScience* + GigaDB (2012)



http://gigasciencejournal.com

http://gigadb.org/

# Lessons Learned

- Technology is really the bottleneck

- Traditional publishing process is too slow & expensive

- Much too focused on narrative and static "version of record"

# A New Approach

**Follow the Software Paradigm?**

CODE   RELEASE   FORK   UPDATE   REPEAT

**Deconstruct the "Version of Record"?**

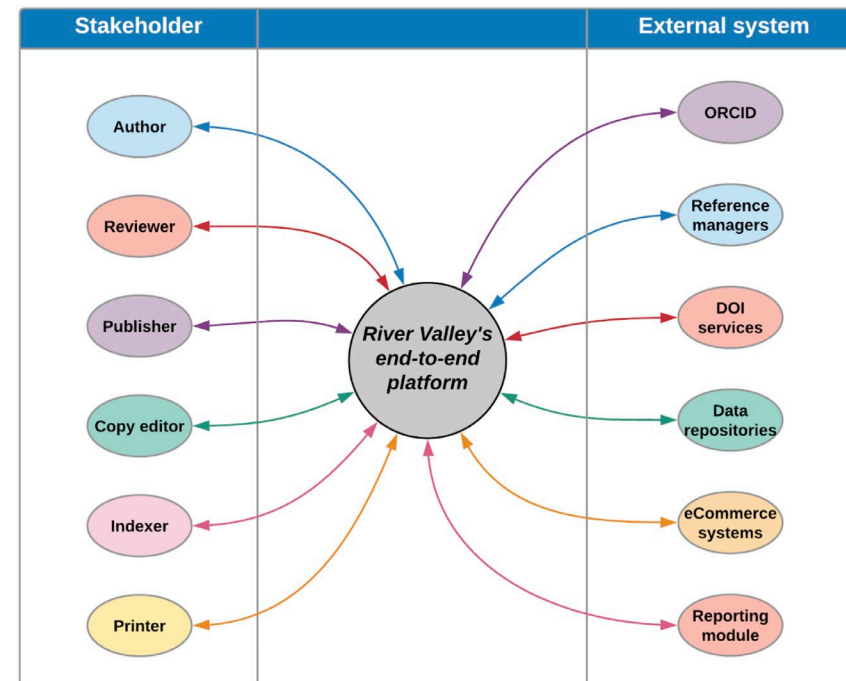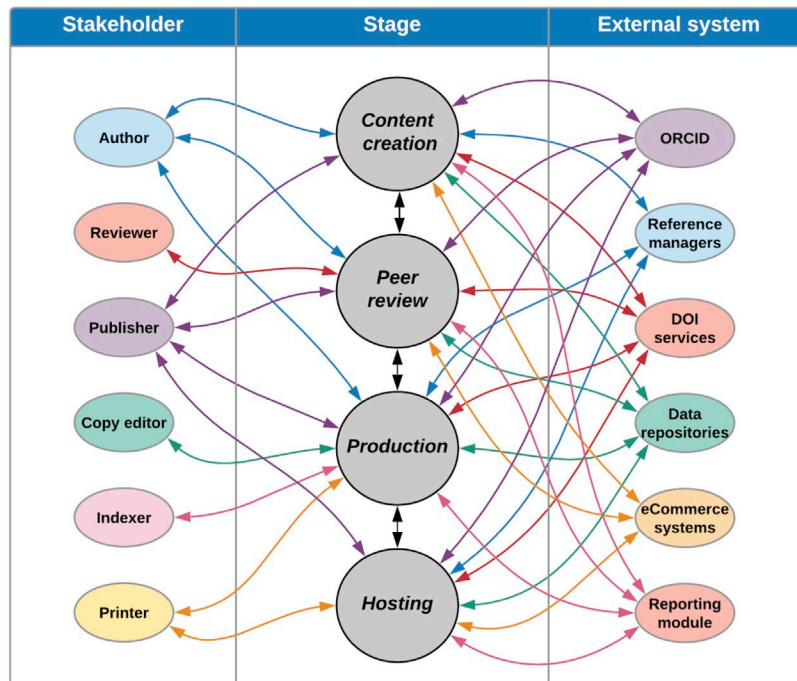DATA   CODE   ENTITIES   FACTS   STABILITY

# Move To New XML End-to-End Pipeline



Custom end-to-end workflow makes integrations simpler with one integration point

# (GIGA)bYte  A New Approach With New Tech

## Main advantage of workflow is XML from start to end

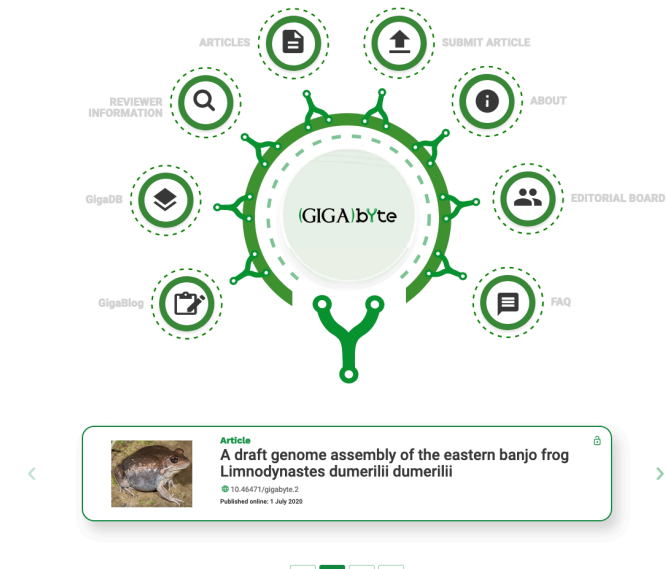Several modules acting as one platform: no import/export of files, so fast and accurate

Cutting out production allows huge time & cost saving (currently 4-8 hours NOT 2 months+)

Any number of versions can be published instantly, including typographic quality PDF or updates/forks –minimal effort

Allows instantaneous switch of views – push of a button

Initial focus on forkable products: data + software + updates

Leverage embeddable dynamic content/widgets

https://gigabytejournal.com/

# Thinking About Users: Authors, Reviewers, Readers

## Reconfigured for short, easy to write & review data & software papers

Streamlined questionnaire-based review

Export as PDF, XML, HTML… "on the fly"



https://gigabytejournal.com/

# Focusing Beyond Version of Record:
# Allows Different Views

# Interactive Features – Increases Understanding and Trust

# Thanks to:

**FOLLOW US:**

**@GigaByteJournal**

**facebook.com/GigaByteJournal**

**http://gigasciencejournal.com/blog/**

editorial@gigabytejournal.com

**Submit now, free APCs till 28th Feb 2021**

**https://gigabytejournal.com/**