# Read Between the Lines:
# Closing Gaps of Materials and Methods to Build Workflow from the Publication

**Tazro Ohta[1], Osamu Ogasawara[2], Yoshinobu Masatani[3], Shigetoshi Yokoyama[3], Kento Aida[3]**

1. Database Center for Life Science, Chiba, Japan 2. DNA DataBank of Japan, Shizuoka, Japan, 3. National Institute of Informatics, Tokyo, Japan

**DBCLS** Database Center for Life Science

**DDBJ** DNA Data Bank of Japan

**NII** 大学共同利用機関法人 情報・システム研究機構 国立情報学研究所 National Institute of Informatics

# Background
## Increasing cost of translation from text to executable scripts

Publishing and sharing data analysis workflow using the galaxy platform has spectacularly reduced the cost of reproducing one's research, but following the description of data analysis which had been performed by other researchers to get the exact same result is still a big challenge. To evaluate the cost of data analysis workflow from the natural language description, we have performed to rebuild the workflow of CAGE sequencing data processing done by FANTOM5 team on the galaxy platform. Though the project has already published a set of papers with a lot of supplementary of methods and online protocols, it was not that straightforward to get the same result from the raw sequencing data available in the public data repository. The results processed by the rebuilt workflow are compared with the results published online by FANTOM5 team. This case study showed that some of the important information to rebuild the workflow is missing even in the well-described documents, for example, the location of the older source code, or the parameters for command execution. As the speed of biological data production increases, it will be more important to build the framework of cost-effective research reproducibility such as an automated evaluation process of published workflow. We will provide the details of our case study, and discuss how we can assure the reproducibility with the galaxy and other possible ways to perform, share, and publish the workflow as it is "executable materials and methods".



Fig. 1. **FANTOM5 (http://fantom.gsc.riken.jp/5/) data processing protocol descriptions.** (a) Description of data processing in the main paper of FANTOM 5 project. (b) Online protocol published by FANTOM 5 project. Though this is helpful, it is still not enough to reproduce the exact same results. (c) A script we made with a help by members of FANTOM 5 project. There are some missing informations from materials and methods or online protocols.
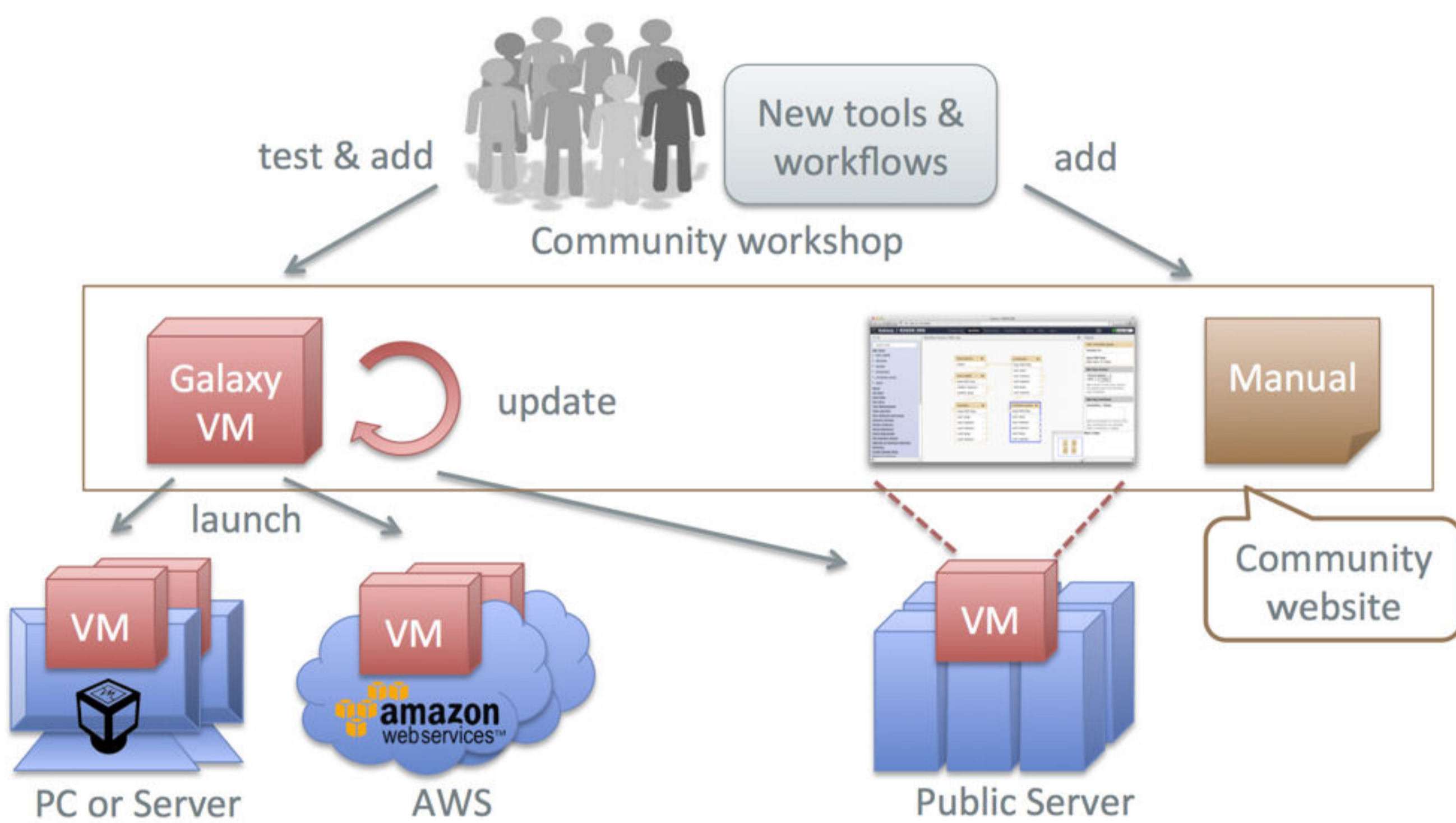


Fig. 2. **Community Galaxy VM and Galaxty Community Japan.** We have been organizing Japanese local community of the galaxy developers and users. The community is sharing the tools and workflows that are used in their labs. Workflows are implemented and fully documented online with the test data. Automated test of the workflow is being developed. Packed VM is distributed as a normal Virtualbox image and Amazon Machine Image via Amazon Web Service. The community is also hosting public server (listed on the page of galaxy project 'virtual appliance').

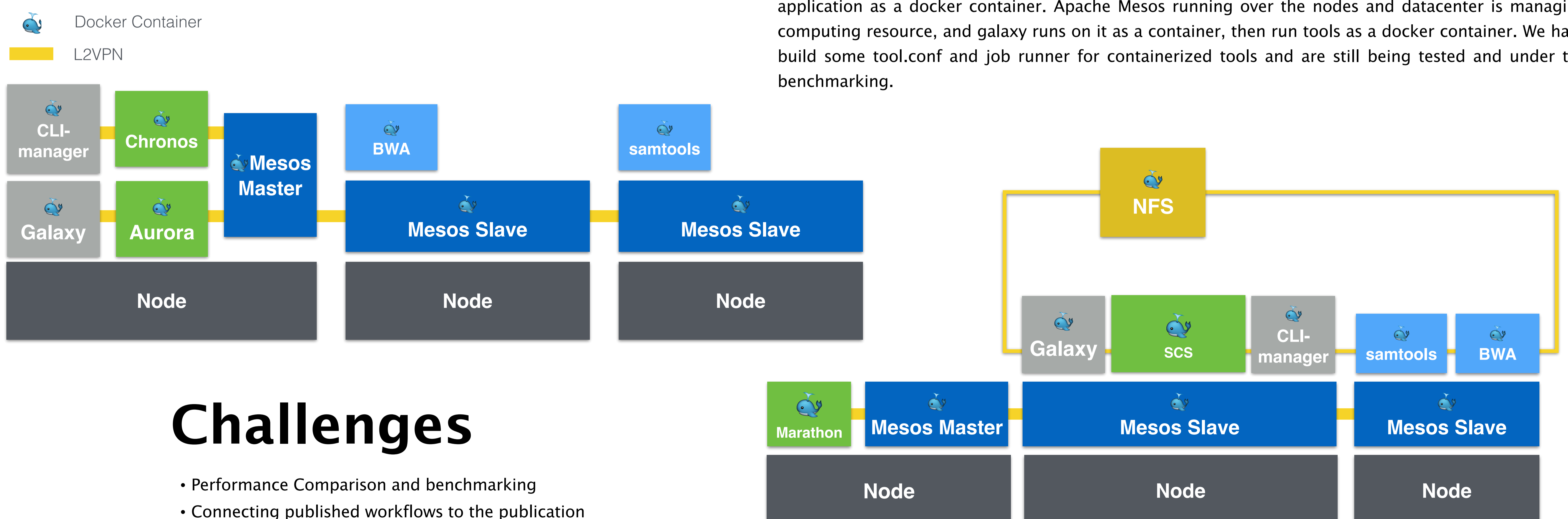# Galaxy as an inter-laboratory workflow sharing platform
## workflow runs on any galaxy, and the galaxy runs anywhere

### Application: Community 'Pitagora' Galaxy VM

We have been trying to develop the system that we are able to share community's workflows developed and used in the various studies which deal with NGS data, such as Exome-Seq, RNA-Seq, ChIP-Seq, or Bisulfite-Seq. We developed Virtual appliance that includes workflows and some tools to help execution of the workflows, and the Virtualbox image and Amazon Machine Image are distributed. We also maintain public server for the test use, and documents and test data are hosted on our project website.

### Platform: 'Overlay Cloud' datacenters

"Overlay Cloud" project led by National Institute of Informatics (Tokyo, Japan) is building large scale cloud computing infrastructure that can distribute computing tasks via their high-speed network system "sinet". We are testing to build galaxy environment on their cloud system, which runs all the process or application as a docker container. Apache Mesos running over the nodes and datacenter is managing computing resource, and galaxy runs on it as a container, then run tools as a docker container. We have build some tool.conf and job runner for containerized tools and are still being tested and under the benchmarking.

# Challenges

- Performance Comparison and benchmarking
- Connecting published workflows to the publication
  - Toolshed or other data repository
- Standard description or annotation of workflows
- Data transfer via internet
  - burst buffer like cloud volume?



Fig. 3. **Overlay cloud system and container-based galaxy and galaxy tools.** (top-left) Current system which is fully containerized computing infrastructure for the galaxy. Each application runs as a container and connected by L2VPN. Data are stored in the NFS mounted on each node. (bottom-right) Current idea of ver. 2.0 system. Apache mesos and marathon container are detached from the network of the containers of galaxy and galaxy tools which are attached to the NFS via L2VPN. This enables more portability to the running system.