



Tazro Ohta

DBCLS, ROIS

twitter.com/inutano

github.com/inutano

speakerdeck.com/inutano



A system to validate published data and studies

Tazro Ohta, Ryota Yamanaka, Osamu Ogasawara, Yoshinobu Masatani, Shigetoshi Yokoyama, Kento Aida

*Reproduction of results from a published study or Meta-analysis using public data is **painful***

What happened during the data analysis

3. Data Processing of Heliscope CAGE data

Sequenced Heliscope reads have a high sequencing error rate ($\sim 5\%$), vary in length and lack an estimation of base qualities. Combined these factors make the data processing challenging. As an initial step we removed reads corresponding to ribosomal RNA. We accomplish this by directly aligning each read against the whole human (mouse) ribosomal DNA complete repeating unit and discarding all reads with an edit distance smaller or equal to two. For this purpose we implemented Myers' bit parallel dynamic programming algorithm¹³ in the program rRNA dust (author: T. Lassmann). For computational efficiency we further parallelized this algorithm using both SIMD instructions and threads. All remaining CAGE reads were mapped to the genome (hg19 and mm9) using Delve, a probabilistic mapper¹⁴. In brief, Delve uses a pair hidden Markov model to iteratively map reads to the genome and estimate position dependent error probabilities. After all error probabilities are estimated, individual reads are placed to a single position on the genome where the alignment has the highest probability to be true according to the pHMM model. Phred scaled mapping qualities¹⁵, reflecting the likelihood of the alignment at a given genome position, are also reported. Reads mapping with a quality of less than 20 (<99% chance of true) were discarded. Furthermore, we discarded all reads that map to the genome with a sequence identity of less than 85%.

*What **exactly** happened during the data analysis*

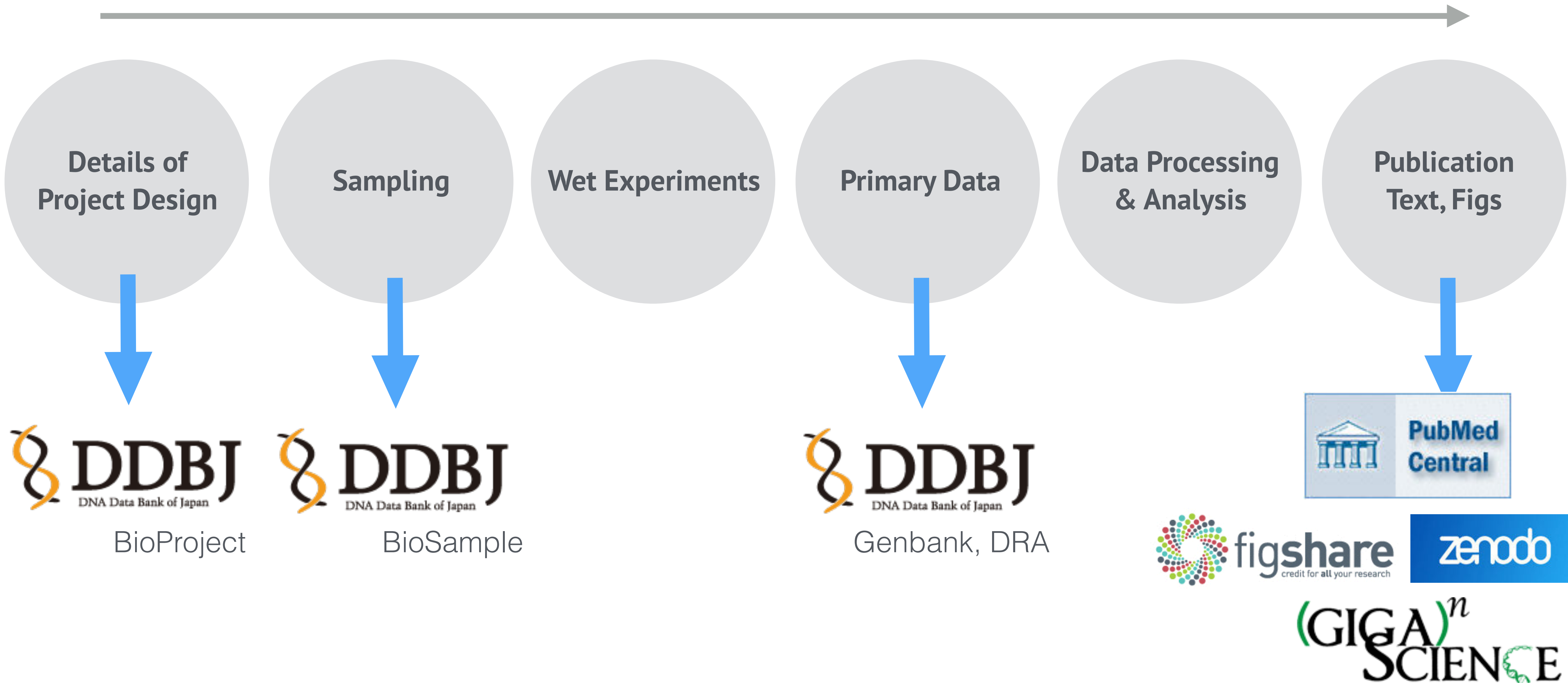
```
74 wget ftp://ftp.ddbj.nig.ac.jp/ddbj_database/dra/fastq/DRA000/DRA000991/DRX008610/DRR009482.fastq.bz
75 bunzip2 DRR009482.fastq.bz2
76
77 head -300000 DRR009482.fastq > input.fastq
78
79 input=input.fastq
80 output=output.sam
81 genome=hg19.fa
82 ribosomal=human_rDNA_U13369.1.fa
83
84 rRNAcust=rRNAfilter/rRNAcust
85 delve=delve/src/delve
86 samtools=samtools-1.2/samtools
87
88 # Cleaning
89 $rRNAcust -s $ribosomal $input -e 2 > tmp1
90
91 # Alignment & Mapping (FASTQ to SAM)
92 $delve index $genome
93 $delve seed tmp1 $genome -o tmp2 -t 8 -l 12 -s 8
94 $delve align tmp2 $genome -u 1 -o $output -t 8
95
96 # SAM to BED
97 $samtools view -bS $output > $output\_bam
```

Missing description:

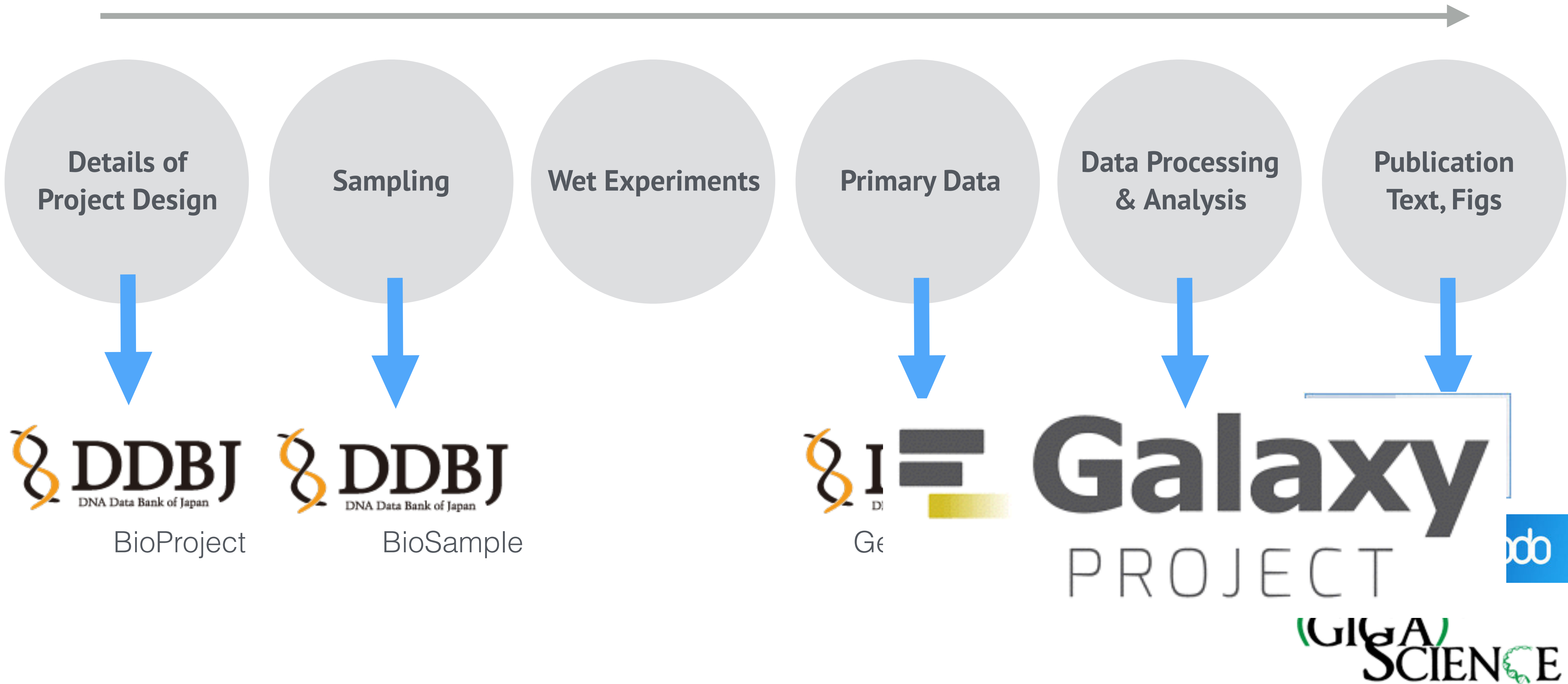
- ***OS/Library dependencies***
- ***Versions of tools, references, databases***
- ***Arguments of tool executions***
- ***Tiny scripts or one liner for data format conversion***

Required: executable materials and methods

Research Activity Time Course



Research Activity Time Course



Research Activity Time Course



Details of Project Design

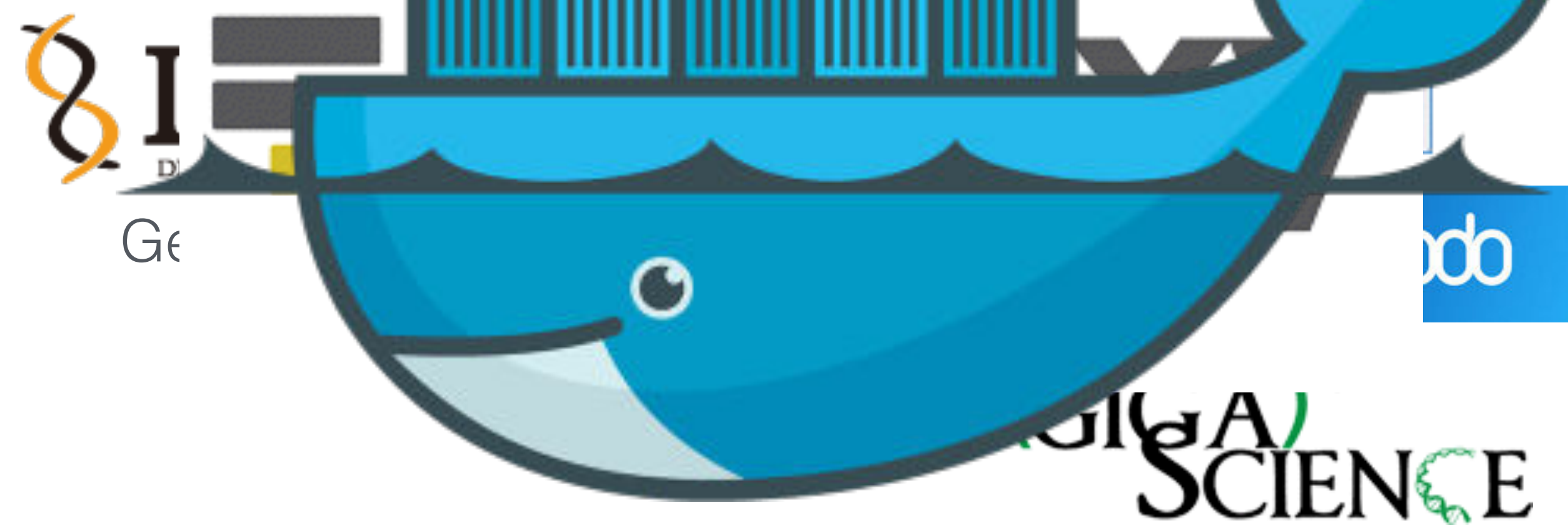
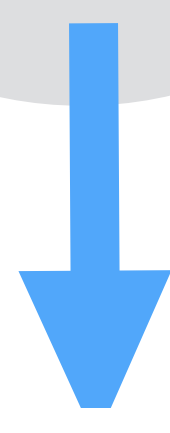
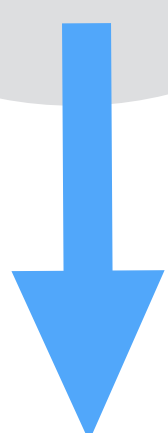
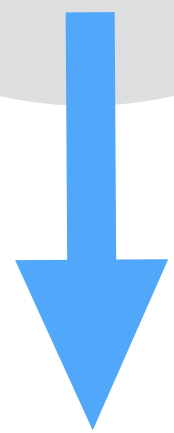
Sampling

Wet Experiments

Primary Data

Data Processing & Analysis

Publication Text, Figs





*proof-of-concept development:
portable, scalable infrastructure for galaxy container*

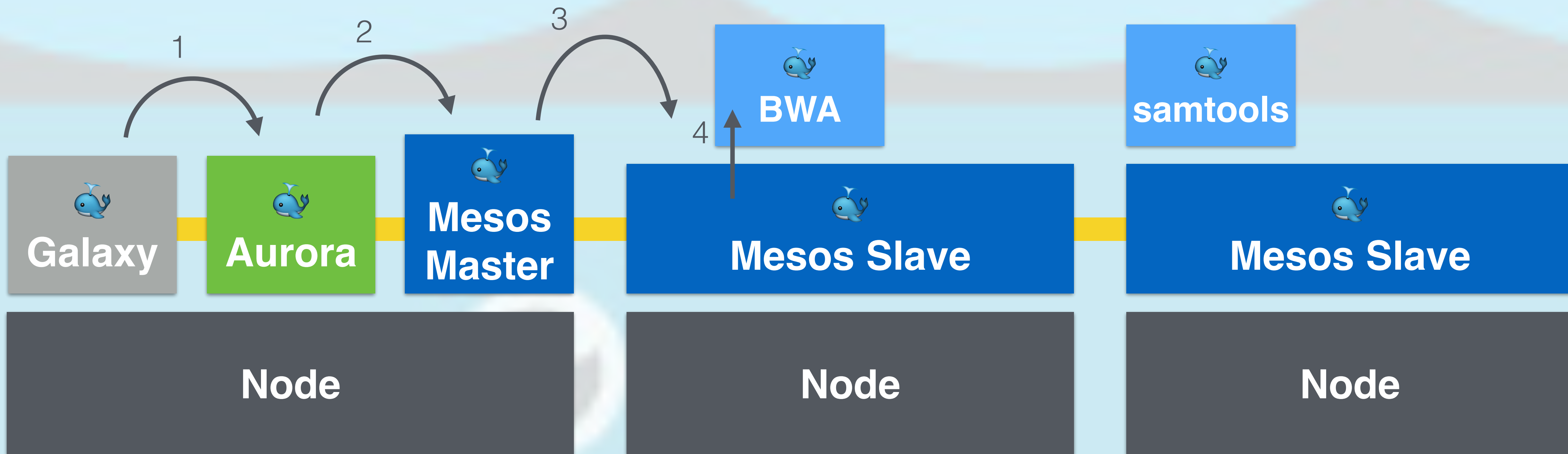


Docker Container



L2VPN

How tools are executed



How data fetched & stored

 Docker Container

 L2VPN

