



The 12th Korea-China-Japan Bioinformatics Training Course & Symposium

Ocean Suites Hotel, Jeju, South Korea
June 18 – 20, 2014

KCJ Bioinformatics Training Course 2014

**”How to analyze Big Data:
Marine Metagenomics and
the Diversity of Microorganisms”**

Ocean Suites Hotel,

June 20, 2014

Takashi Gojobori

**KAUST (King Abdullah University of
Science and Technology), KSA and
NIG (National Institute of Genetics), Japan**

How to analyze Big Data (1)!

**Genome analysis is
a tool,
but not purpose!**

NGS

Bio-samples

**Cells,
Tissues,
Organs
Species
Popula-
tions**

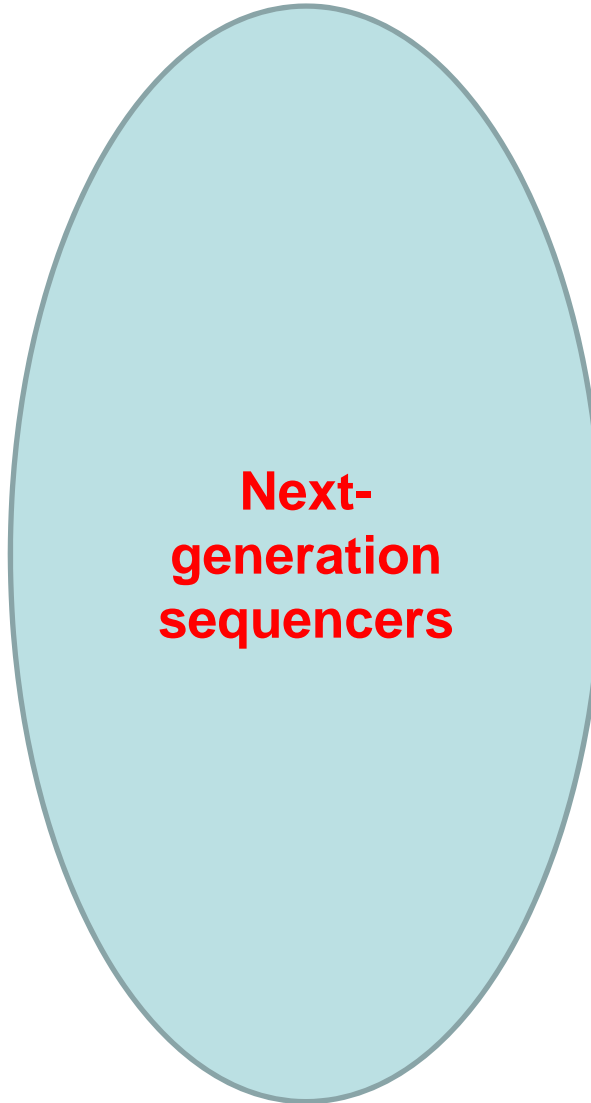
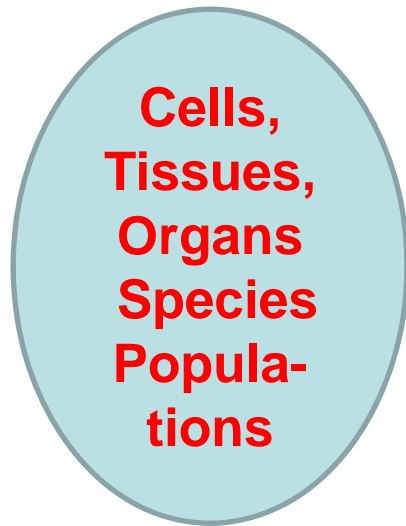


**Next-
generation
sequencers**



**Database
Data Analysis
Informatics**

**Seque
nce
data**



Item	Description
Read Length and Speed	512 nanopores x 15bp/sec => ~7500 bp/sec
Read Accuracy	99.8%
6 Hours Life Time	150 x 106bp
Applied Current / Blockage	60 picoamps to anywhere from 20-40 picoamps
No. of nanopore	2,000 nanopores / cartridge. Will become available in early 2013 containing over 8,000 nanopores. →Delivers a complete human genome in 15 minutes.
Sample Preparation	Any user-derived sample preparation resulting in double stranded DNA (dsDNA) in solution is compatible with the system.
Amplification	No sample amplification.
Cost	\$900
Commercialization	Oxford Nanopore intends to sell the system directly to customers with



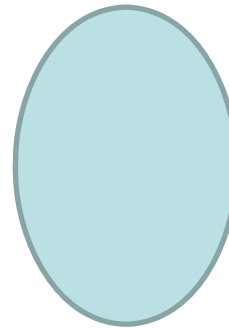
Nano Pore Oxford (2012)

Bio-samples

**/Cells,
/Tissues,
/Organs
/ Species
/Popula-tions
+
/Time
/Environ-
ments
/Conditions**



NGS



Database Data Analysis Informatics

**/Genome
/Meta-
genome
/Epi-genome
/RNA-seq
/CHIP-seq
/PPI
/Synthe-tic**

How to analyze Big Data (2)!

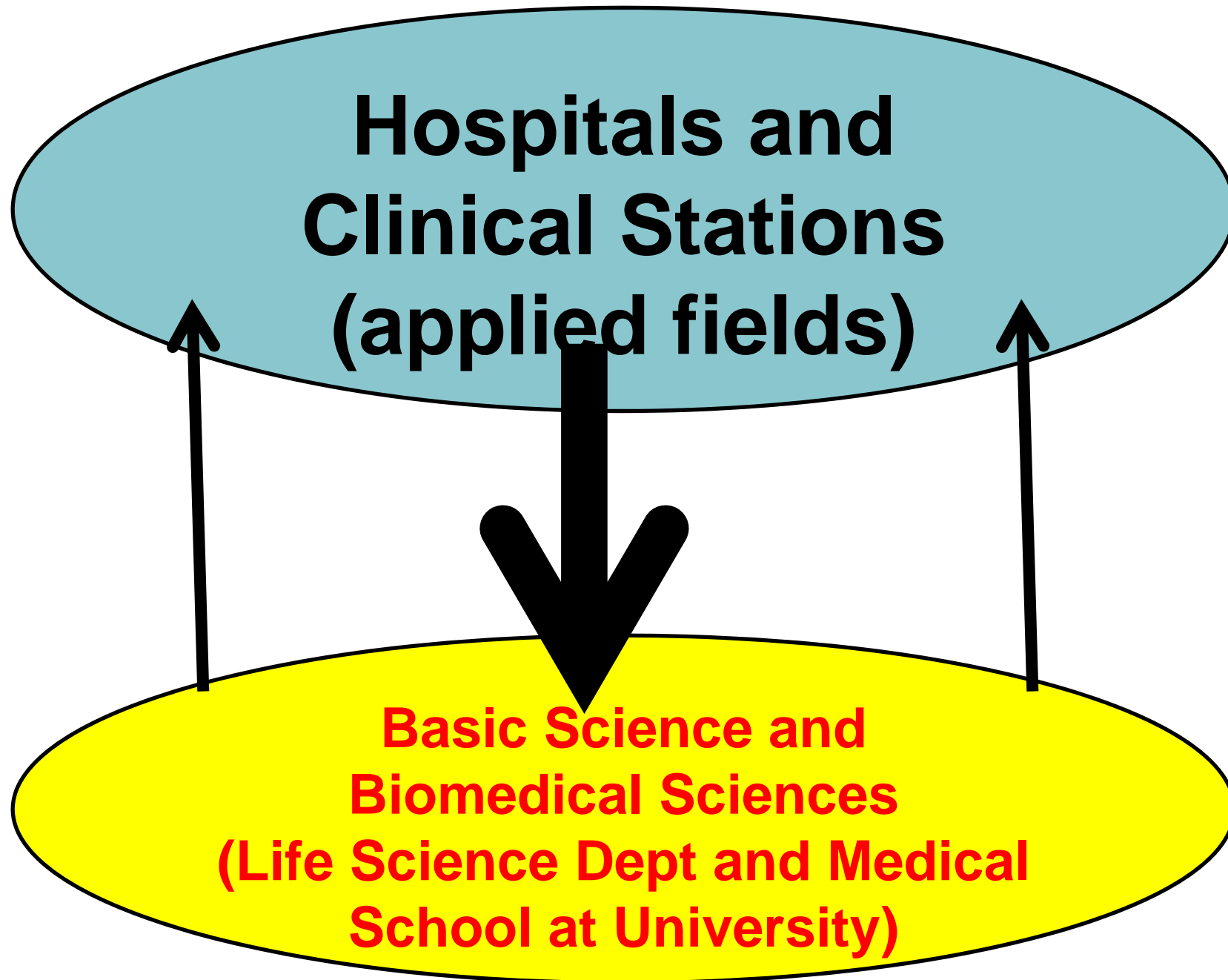
**Biosamples are
a Key!**

Large or rare

**Hospitals and
Clinical Stations
(applied fields)**

The diagram consists of two ovals. The top oval is light blue and contains the text 'Hospitals and Clinical Stations (applied fields)'. The bottom oval is yellow and contains the text 'Basic Science and Biomedical Sciences (Life Science Dept and Medical School at University)'. Two black arrows point upwards from the yellow oval to the blue oval, indicating a flow or relationship from basic science to clinical applications.

**Basic Science and
Biomedical Sciences
(Life Science Dept and Medical
School at University)**

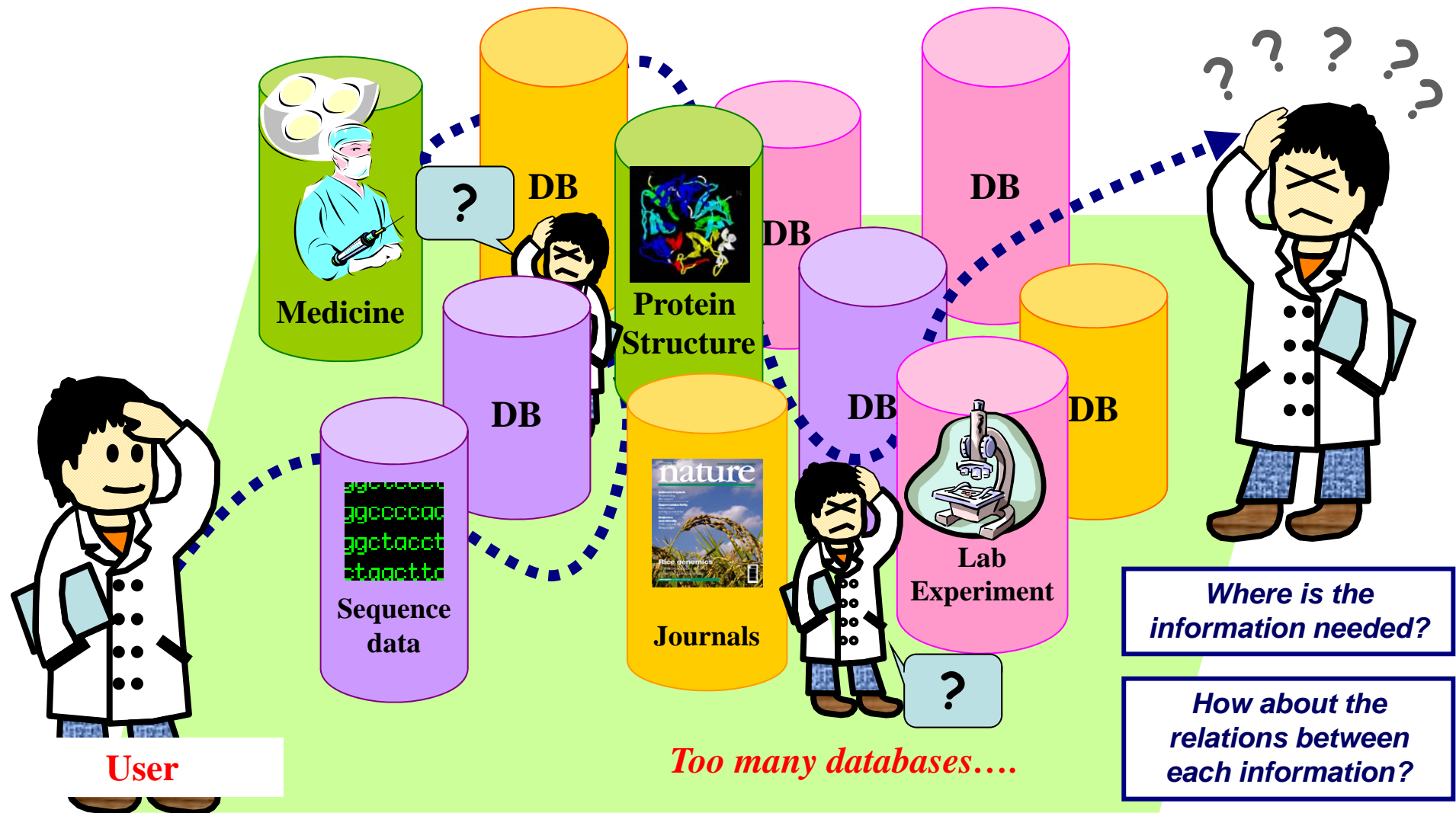


How to analyze Big Data (3)!

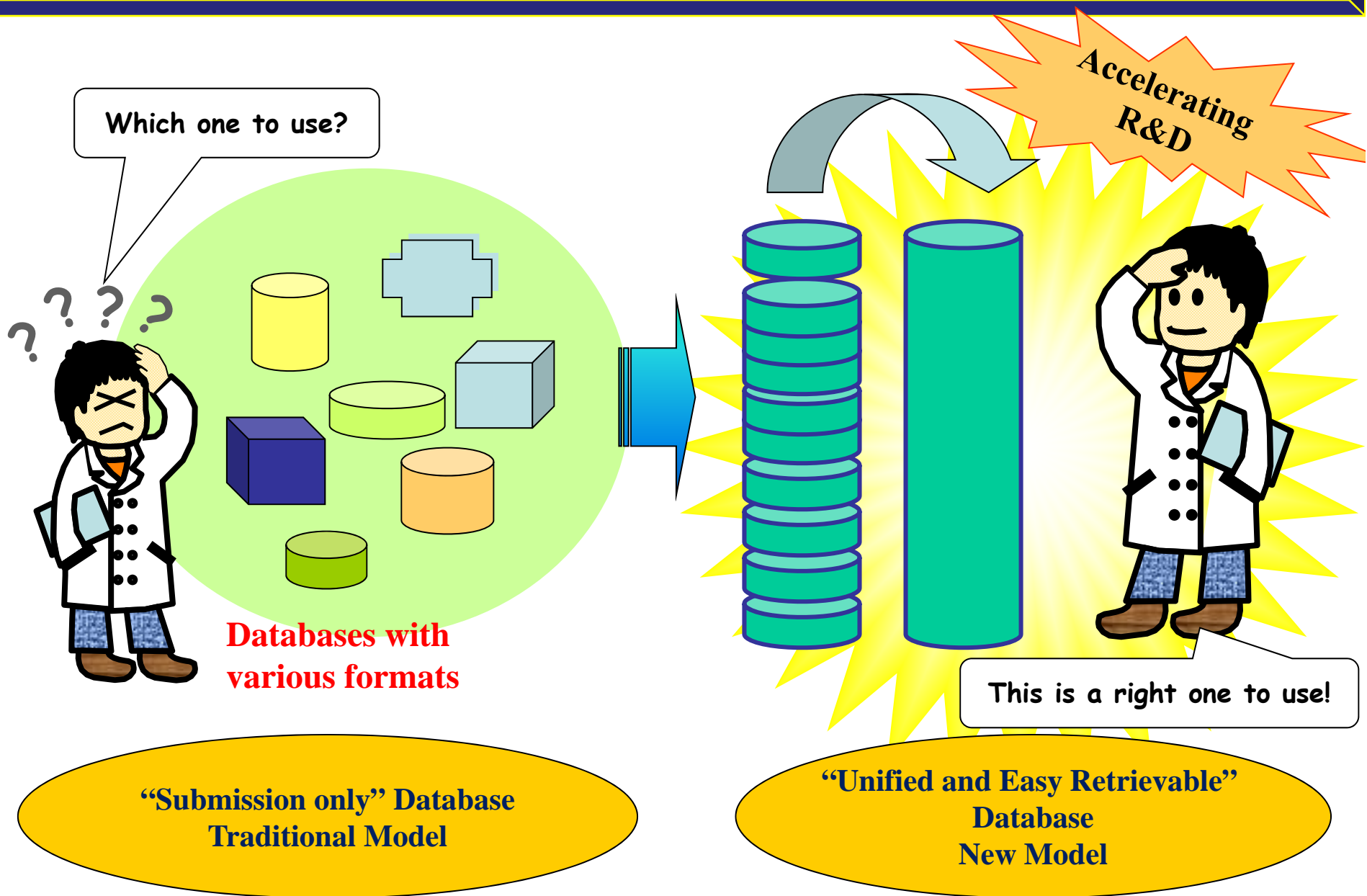
**Database
Construction is
a Key!**

Issues on retrieving the necessary information

(The Gist) **Lack of the standard format without unified information often hinders research and development seriously**



Be Unified and Easily Retrievable Format from Sporadic Information !





D. Howe, M. Costanzo, P. Fey, **T. Gojobori,
L. Hannick, W. Hide, D. Hill, R. Kania, M. Schaeffer,
S. St Pierre, S. Tweigger, and S. Rhee**

Nature (2008) 455: 47-50

How to analyze Big Data (4)!

**A gene for the genome is
a Key!**

Evolutionary changes of multiple visual pigment genes in the complete genome of Pacific bluefin tuna

Yoji Nakamura^{a,1}, Kazuki Mori^b, Kenji Saitoh^a, Kenshiro Oshima^c, Miyuki Mekuchi^a, Takuma Sugaya^a, Yuya Shigenobu^a, Nobuhiko Ojima^a, Shigeru Muta^b, Atushi Fujiwara^a, Motoshige Yasuike^a, Ichiro Oohara^a, Hideki Hirakawa^d, Vishwajit Sur Chowdhury^e, Takanori Kobayashi^f, Kazuhiro Nakajima^g, Motohiko Sano^a, Tokio Wada^f, Kousuke Tashiro^b, Kazuho Ikee^h, Masahira Hattori^c, Satoru Kuhara^b, Takashi Gojobori^{h,1}, and Kiyoshi Inouye^f

^aNational Research Institute of Fisheries Science, Fisheries Research Agency, Kanazawa, Yokohama 236-8648, Japan; ^bDepartment of Bioscience and Biotechnology, Faculty of Agriculture, Kyushu University, Higashi-ku, Fukuoka 812-8581, Japan; ^cGraduate School of Frontier Sciences, University of Tokyo, Kashiwa 277-8581, Japan; ^dDepartment of Plant Genome Research, Kazusa DNA Research Institute, Kisarazu 292-0818, Japan; ^eInternational Education Center, Faculty of Agriculture, Kyushu University, Higashi-ku, Fukuoka 812-8581, Japan; ^fFisheries Research Agency, Nishi, Yokohama 220-8113, Japan; ^gJapan Sea National Fisheries Research Institute, Chuou-ku, Niigata 951-8121, Japan; and ^hCenter for Information Biology, National Institute of Genetics, Mishima 411-8540, Japan

Edited* by Tomoko Ohta, National Institute of Genetics, Mishima, Japan, and approved May 20, 2013 (received for review February 2, 2013)

Tunas are migratory fishes in offshore habitats and top predators with unique features. Despite their ecological importance and high market values, the open-ocean lifestyle of tuna, in which effective sensing systems such as color vision are required for capture of prey, has been poorly understood. To elucidate the genetic and evolutionary basis of optic adaptation of tuna, we determined the genome sequence of the Pacific bluefin tuna (*Thunnus orientalis*), using next-generation sequencing technology. A total of 26,433 protein-coding genes were predicted from 16,802 assembled scaffolds. From these, we identified five common fish visual pigment genes: red-sensitive (middle/long-wavelength sensitive; M/LWS), UV-sensitive (short-wavelength sensitive 1; SWS1), blue-sensitive (SWS2), rhodopsin (RH1), and green-sensitive (RH2) opsin genes. Sequence comparison revealed that tuna's RH1 gene has an amino acid substitution that causes a short-wave shift in the absorption spectrum (i.e., blue shift). Pacific bluefin tuna has at least five RH2 paralogs, the most among studied fishes; four of the proteins encoded may be tuned to blue light at the amino acid level. Moreover, phylogenetic analysis suggested that gene conversions have occurred in each of the SWS2 and RH2 loci in a short period. Thus, Pacific bluefin tuna has undergone evolutionary changes in three genes (RH1, RH2, and SWS2), which may have contributed to detecting blue-green contrast and measuring the distance to prey in the blue-pelagic ocean. These findings provide basic information on behavioral traits of predatory fish and, thereby, could help to improve the technology to culture such fish in captivity for resource management.

tuna genome

Tunas are considered “the ultimate fish,” because they are top predators in ocean ecosystems, in addition to their unique

duplication occurred and each copy has been maintained for a long time. Previous studies have accumulated molecular information on opsins, focusing on the residues surrounding the retinal-binding pocket (10–12).

Many physiological studies have been conducted on fish visual system underwater (13, 14). As for tuna, spectrophotometric analyses have demonstrated several wavelengths of maximal absorbance (λ_{max}) in the visual pigments of yellowfin tuna (*Thunnus albacares*) (15) and Pacific bluefin tuna (*Thunnus orientalis*) (16), respectively. However, little is known about the genetic basis and evolutionary history of tuna's optic adaptation to an open-ocean predatory lifestyle. In this study, we have sequenced the draft genome of Pacific bluefin tuna and analyzed the opsin genes in a phylogenetic framework to look for evidence of optic adaptation at the molecular level. The origins of tuna's opsin paralogs were dated by genome-wide comparison among the teleosts for which genomic information is available, and the relationship between the evolutionary pathway of opsin genes and adaptation to ocean environment is discussed.

Results

Genome Sequencing and Gene Prediction. The diploid genome of a wild-caught male Pacific bluefin tuna (*T. orientalis*) was sequenced. A whole-genome shotgun sequencing and assembling strategy with a combination of Roche 454 FLX Titanium (Roche Diagnostics) and Illumina Gallx platforms provided 192,169 contigs (>500 bp) and 16,802 scaffolds (>2 kb) totaling 740.3 Mb, corresponding to 92.5% of the estimated genome size (~800 Mb) (ref. 17; Table 1). The scaffolds were obtained by assembling 31.9 million 454 reads, including 4.9 million paired ends (11.9-fold coverage) and 229.7 million Illumina paired-end reads (43-fold coverage) (Tables S1 and S2 and Fig. S1). Sequence

How to analyze Big Data (5)!

**Comparative (medical,
evolutionary) genome is
a Key!**

The First Symbiont-Free Genome Sequence of Marine Red Alga, Susabi-nori (*Pyropia yezoensis*)

Yoji Nakamura^{1*}, Naobumi Sasaki², Masahiro Kobayashi³, Nobuhiko Ojima¹, Motoshige Yasuike¹, Yuya Shigenobu¹, Masataka Satomi¹, Yoshiya Fukuma⁴, Koji Shiwaku⁴, Atsumi Tsujimoto⁵, Takanori Kobayashi⁶, Ichiro Nakayama⁷, Fuminari Ito⁸, Kazuhiro Nakajima⁹, Motohiko Sano¹, Tokio Wada⁶, Satoru Kuhara¹⁰, Kiyoshi Inouye⁶, Takashi Gojobori^{2*}, Kazuho Ikeo²

1 National Research Institute of Fisheries Science, Fisheries Research Agency, Yokohama, Kanagawa, Japan, **2** Center for Information Biology, National Institute of Genetics, Mishima, Shizuoka, Japan, **3** Seto National Fisheries Research Institute, Fisheries Research Agency, Nagasaki-shi, Nagasaki, Japan, **4** Hitachi Solutions, Ltd., Shinagawa-ku, Tokyo, Japan, **5** Japan Software Management Co. Ltd., Yokohama, Kanagawa, Japan, **6** Fisheries Research Agency, Yokohama, Kanagawa, Japan, **7** Ministry of Agriculture, Forestry and Fisheries, Chiyoda-ku, Tokyo, Japan, **8** National Research Institute of Aquaculture, Fisheries Research Agency, Minami-Ise, Mie, Japan, **9** Japan Sea National Fisheries Research Institute, Fisheries Research Agency, Chuou-ku, Niigata, Japan, **10** Division of Molecular Biosciences, Department of Bioscience and Biotechnology, Faculty of Agriculture, Kyushu University, Higashi-ku, Fukuoka, Japan

Abstract

Nori, a marine red alga, is one of the most profitable mariculture crops in the world. However, the biological properties of this macroalga are poorly understood at the molecular level. In this study, we determined the draft genome sequence of susabi-nori (*Pyropia yezoensis*) using next-generation sequencing platforms. For sequencing, thalli of *P. yezoensis* were washed to remove bacteria attached on the cell surface and enzymatically prepared as purified protoplasts. The assembled contig size of the *P. yezoensis* nuclear genome was approximately 43 megabases (Mb), which is an order of magnitude smaller than the previously estimated genome size. A total of 10,327 gene models were predicted and about 60% of the genes validated lack introns and the other genes have shorter introns compared to large-genome algae, which is consistent with the compact size of the *P. yezoensis* genome. A sequence homology search showed that 3,611 genes (35%) are functionally unknown and only 2,069 gene groups are in common with those of the unicellular red alga, *Cyanidioschyzon merolae*. As color trait determinants of red algae, light-harvesting genes involved in the phycobilisome were predicted from the *P. yezoensis* nuclear genome. In particular, we found a second homolog of phycobilisome-degradation gene, which is usually chloroplast-encoded, possibly providing a novel target for color fading of susabi-nori in aquaculture. These findings shed light on unexplained features of macroalgal genes and genomes, and suggest that the genome of *P. yezoensis* is a promising model genome of marine red algae.

Citation: Nakamura Y, Sasaki N, Kobayashi M, Ojima N, Yasuike M, et al. (2013) The First Symbiont-Free Genome Sequence of Marine Red Alga, Susabi-nori (*Pyropia yezoensis*). PLoS ONE 8(3): e57122. doi:10.1371/journal.pone.0057122

Editors: Juergen Kroymann, French National Centre for Scientific Research, Université Paris-Sud, France

Received: November 30, 2012; **Accepted:** January 9, 2013; **Published:** March 11, 2013

Copyright: © 2013 Nakamura et al. This is an open-access article distributed under the terms of the Creative Commons Attribution License, which permits unrestricted use, distribution, and reproduction in any medium, provided the original author and source are credited.

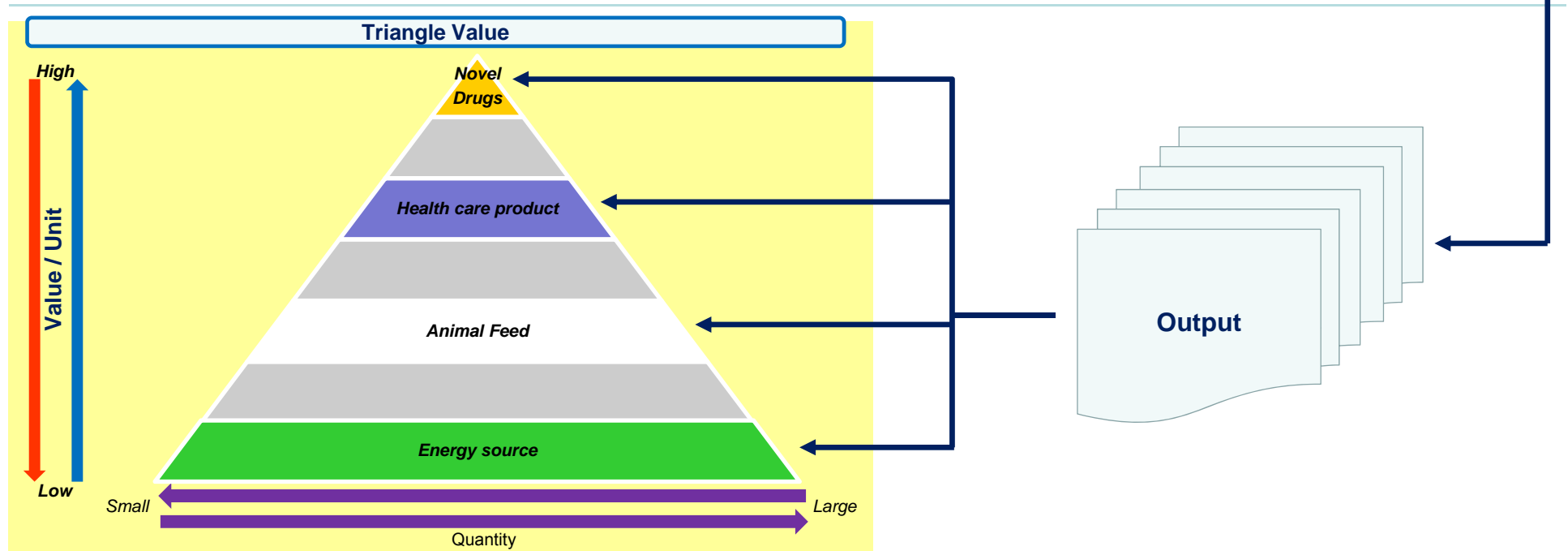
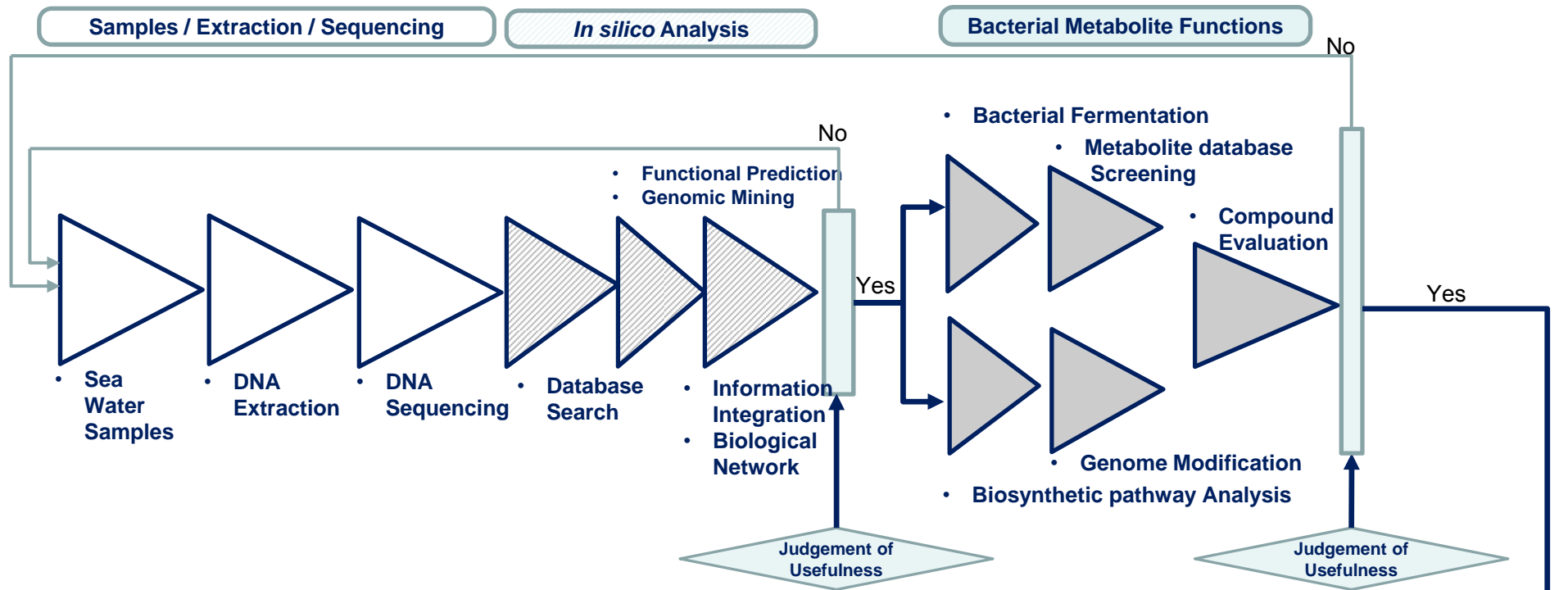
Funding: This work was supported by grants from the Fisheries Agency, Ministry of Agriculture, Forestry and Fisheries, Japan, and from the Fisheries Research Agency, Japan. The funders had no role in study design, data collection and analysis, decision to publish, or preparation of the manuscript.

~Marine Micro Meta- Genomics~

[写真] Apollo 17号からみた地球 (NASA)

How to analyze Big Data (6)!

**A pipeline is
a Key!**

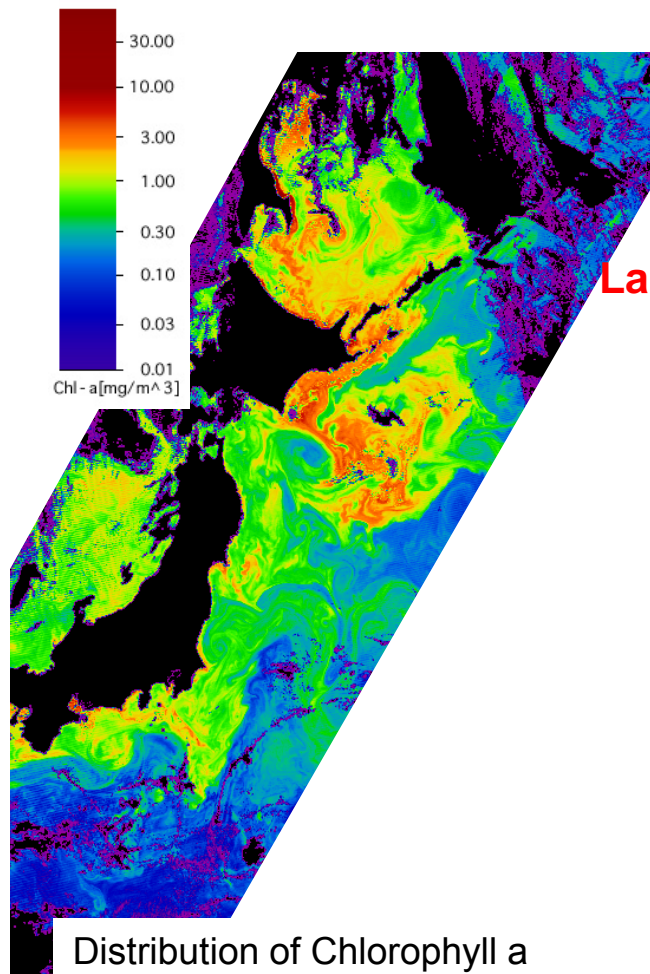


How to analyze Big Data (7)!

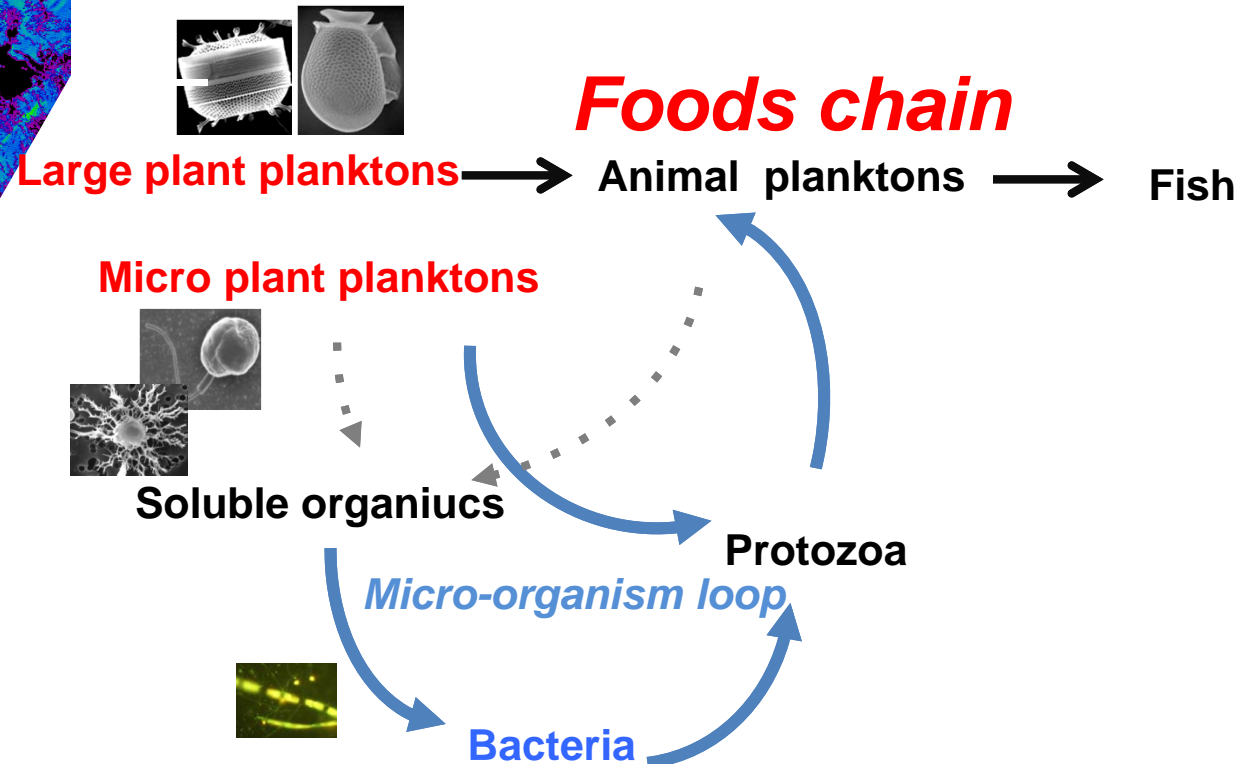
**“Science is gamble,
then need to win” is
a Key!**

—A View of marine micro-organism ecosystem—

Tohoku sea coast → One of species-rich spots



Distribution of Chlorophyll a
In Tohoku sea coast



Understanding of Marine Microorganism Diversity by use of the Digital DNA chip system through Metagenomics

～Subjects and Division of Roles～

Subjects		Role
1	Data Analysis of Metagenomics and Project Management	NIG
2	Construction of Marine Metagenomics Database and Developments of the Digital DNA Chip System	NIG/JSM
3	NGS Sequencing of Metagenomic data	Tokyo U/ NIG
4	Isolation of DNAs from sea water samples/ Sampling machine developments	Kyushu U
5	Sea Water Sampling and Measurements of Physical Conditions	NMRC
6	Analysis of Plant Pico-Planktons by Flow Cytometry	NIE







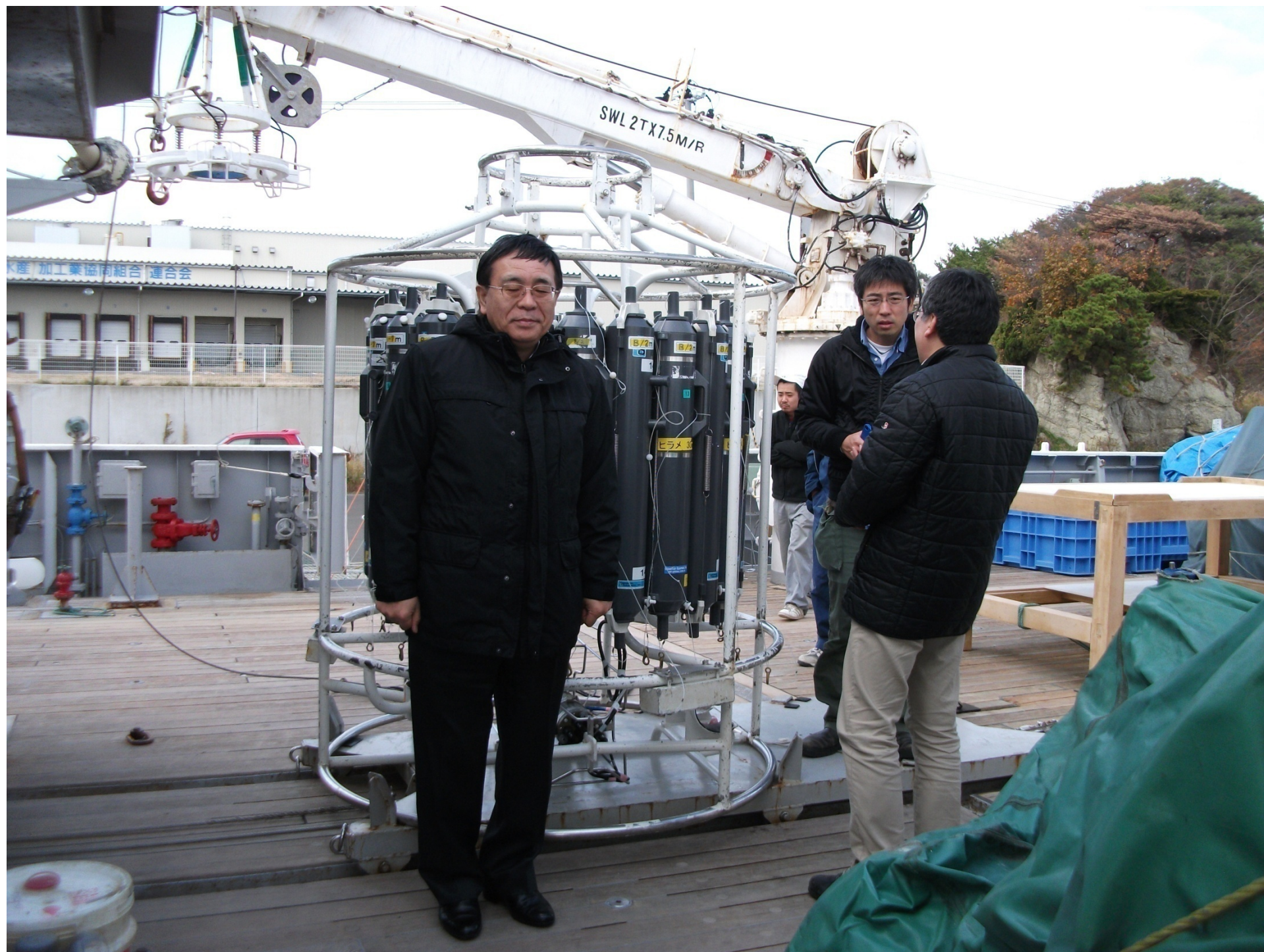
丸 鷹 若

WAKATAKA MARU





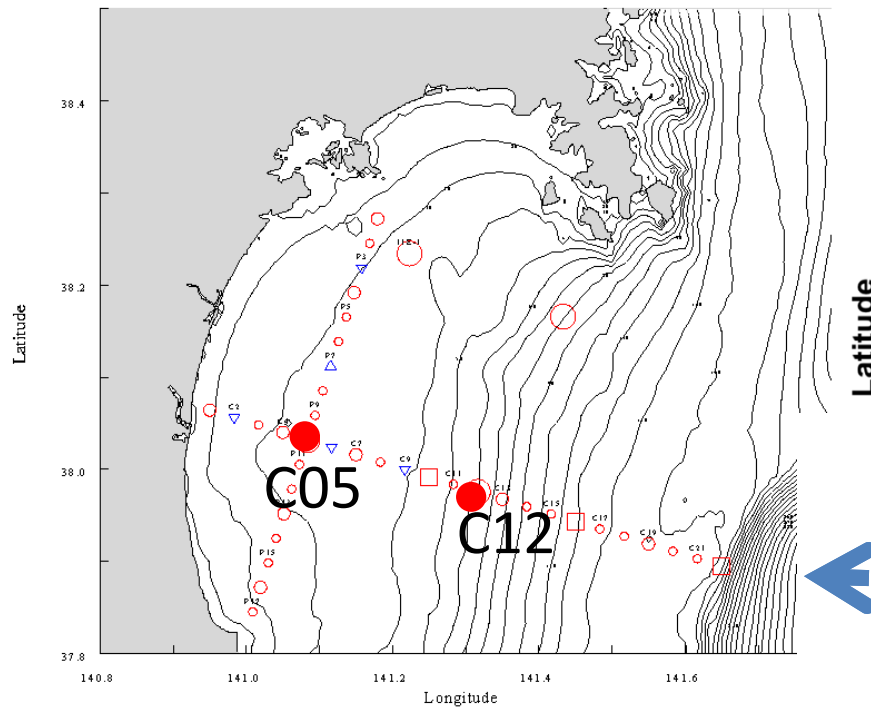




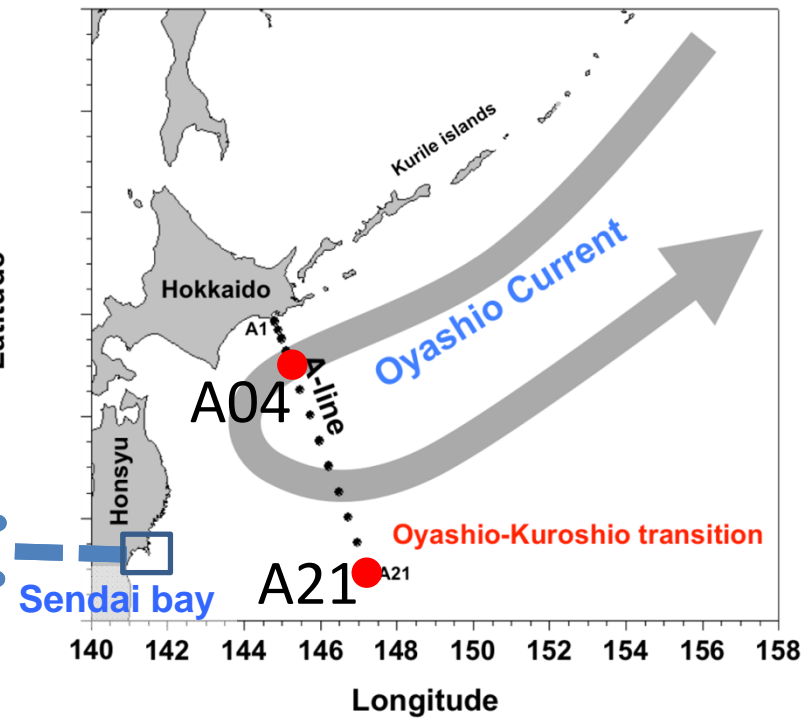




Observation Points at the Sendai Bay and Nemuro Sea as a Control



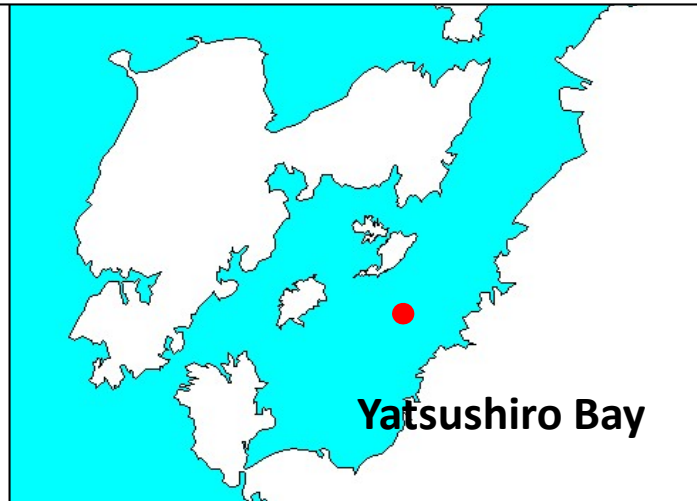
Observation Points (Sendai Bay)



A-line (Near Sea)

Sea Water Sampling Points in Kyushu for “Red Tide” Monitoring

Red Tide Occurrence Sea Area



Red Tide Non-Occurrence Sea Area



Comparative Metagenomics for
Detecting any changes in
marine microorganism diversity

Periodical Sampling

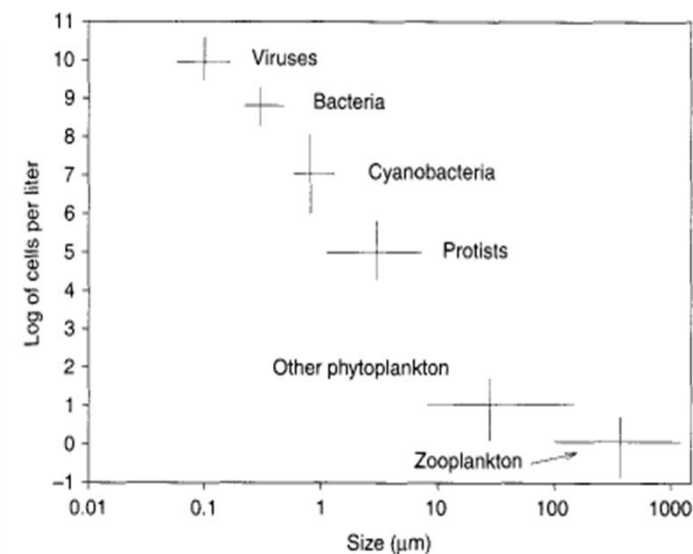


Sampling of Sea Waters at Observation Points in Sendai Bay and Nemuro Sea

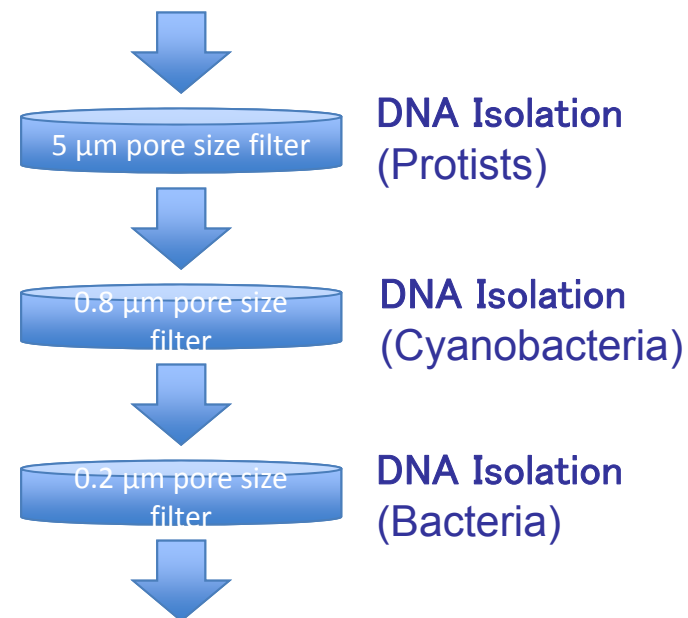
Sails	Dates	Sea points	Depth (m)	Sea Volume (Kg)
WK1203	2012.3.5	A04	10	22.8
WK1203	2012.3.5	A04	30	20.6
WK1203	2012.3.11	A21	10	20.0
WK1203	2012.3.11	A21	10	20.0
WK1203	2012.3.13	P10(C5)	1	7.0
WK1203	2012.3.13	C12	1	21.0
Rent1204	2012.4.16	P10(C5)	1	15.1
Rent1204	2012.4.16	P10(C5)	18	6.8
Rent1204	2012.4.16	C12	1	4.5
Rent1204	2012.4.16	C12	10	7.2
WK1205	2012.5.10	A04	10	17.1
WK1205	2012.5.10	A04	20	15.1
WK1205	2012.5.15	A21	10	6.8
WK1205	2012.5.15	A21	40	9.0
WK1205	2012.5.16	P10(C5)	1	9.2
WK1205	2012.5.16	P10(C5)	15	9.3
WK1205	2012.5.17	C12	1	9.4
WK1205	2012.5.17	C12	20	8.6
WK1206D	2012.6.17	F3	5	8.5
WK1206D	2012.6.17	F3	20	8.0
WK1206D	2012.6.18	P10(C5)	5	9.0
WK1206D	2012.6.18	P10(C5)	20	9.0
WK1207	2012.7.14	C5(P10)	2	8.0
WK1207	2012.7.14	C5(P10)	10	8.5
WK1207	2012.7.16	C12	2	9.0
WK1207	2012.7.16	C12	30	9.5

Isolation of DNAs (Kyushu U.)

Sails	Dates	Points	Depth (m)	Volume (Kg)
WK1203	2012.3.5	A04	10	22.8
WK1203	2012.3.5	A04	30	20.6
WK1203	2012.3.11	A21	10	20.0
WK1203	2012.3.11	A21	10	20.0
WK1203	2012.3.13	P10 (C5)	1	7.0
WK1203	2012.3.13	C12	1	21.0
Rent1204	2012.4.16	P10 (C5)	1	15.1
Rent1204	2012.4.16	P10 (C5)	18	6.8
Rent1204	2012.4.16	C12	1	4.5
Rent1204	2012.4.16	C12	10	7.2
WK1205	2012.5.10	A04	10	17.1
WK1205	2012.5.10	A04	20	15.1
WK1205	2012.5.15	A21	10	6.8
WK1205	2012.5.15	A21	40	9.0
WK1205	2012.5.16	P10 (C5)	1	9.2
WK1205	2012.5.16	P10 (C5)	15	9.3
WK1205	2012.5.17	C12	1	9.4
WK1205	2012.5.17	C12	20	8.6
WK1206D	2012.6.17	F3	5	8.5
WK1206D	2012.6.17	F3	20	8.0
WK1206D	2012.6.18	P10 (C5)	5	9.0
WK1206D	2012.6.18	P10 (C5)	20	9.0
WK1207	2012.7.14	C5 (P10)	2	8.0
WK1207	2012.7.14	C5 (P10)	10	8.5
WK1207	2012.7.16	C12	2	9.0
WK1207	2012.7.16	C12	30	9.5



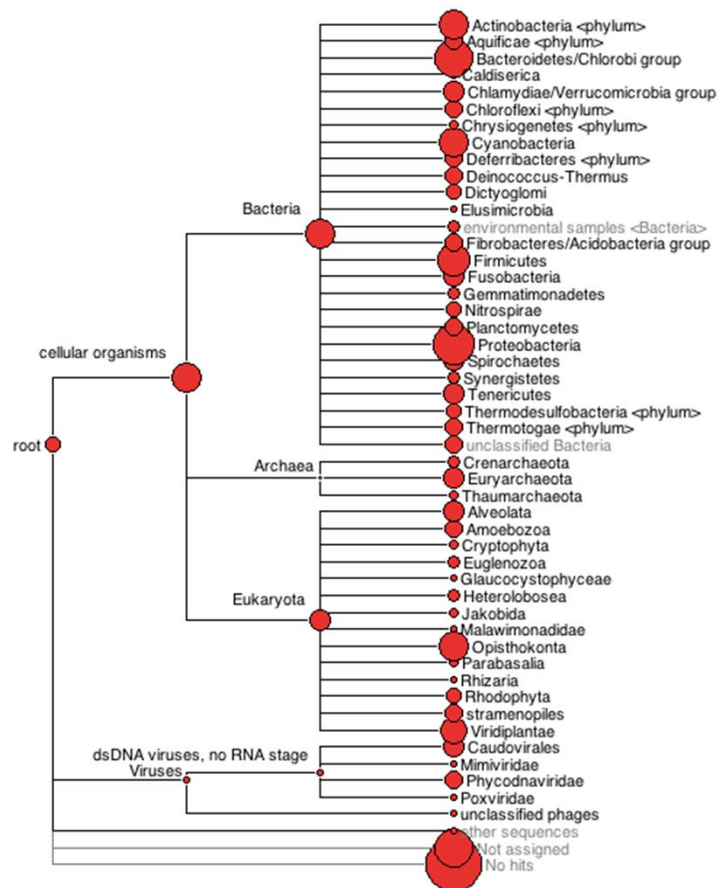
Sea water samples



Metagenomic Data Analysis –1– (NIG)

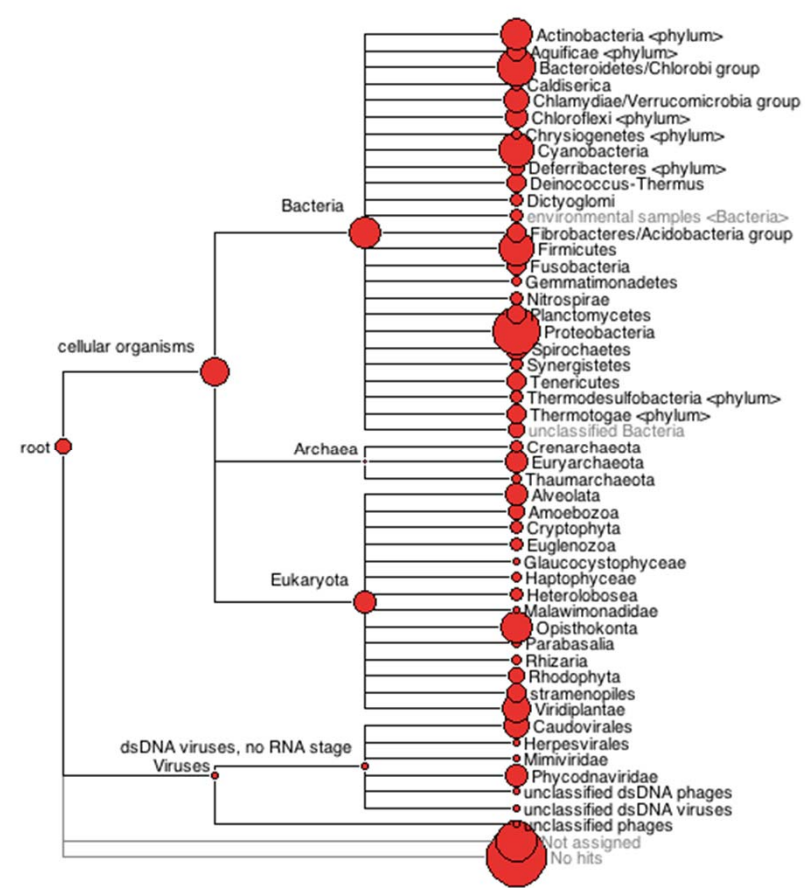
Buzen-1

Total reads: 2899428
Assigned reads: 102789
Unassigned reads: 18583
Reads with no hits: 2778056



Buzen-2

Total reads: 4340338
Assigned reads: 168021
Unassigned reads: 26405
Reads with no hits: 4145912



Conclusion and Summary

Summary

- 1. We constructed pipelines for understanding dynamics of marine bacteria in Sendai area using metagenomic data analysis**
- 2. Proteobacteria was major bacteria in spring and winter while cyanobacteria was increased in August and October in Sendai area**
- 3. Bacterial communities at surface were similar to those at 20m below surface in spring and winter, while bacterial communities at surface in summer and autumn were different from those at 20m below surface**
- 4. Metagenomics showed that 26 metabolic pathways were present in the Sendai area and genes associated with photosynthesis pathway were increased in August and October**

Conclusion

**“Metagenomics is a powerful method
for understanding dynamics of
marine bacteria in ocean ecosystem”**

How to analyze Big Data (8)!

**Monitoring (time and space)
is
a Key!**

Sampling device with a monitoring system (Kyushu U.)

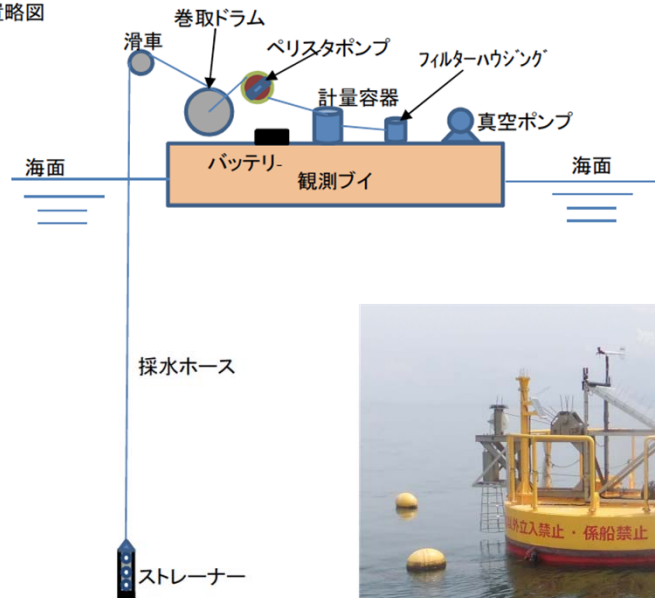
ウインチ装置部参考図面

自立式海洋微生物DNA採取装置

装置製作条件

表層から10mまでの柱状サンプル(2.5L程度)を1日に1回採取し、既定のフィルターでろ過を行う。
これを1週間行った後、サンプルとしてフィルター(7個)を回収する。
電源に関しては、今回は太陽電池等を用いないため動作に必要な電源を、鉛蓄電池で確保する。

採取装置略図

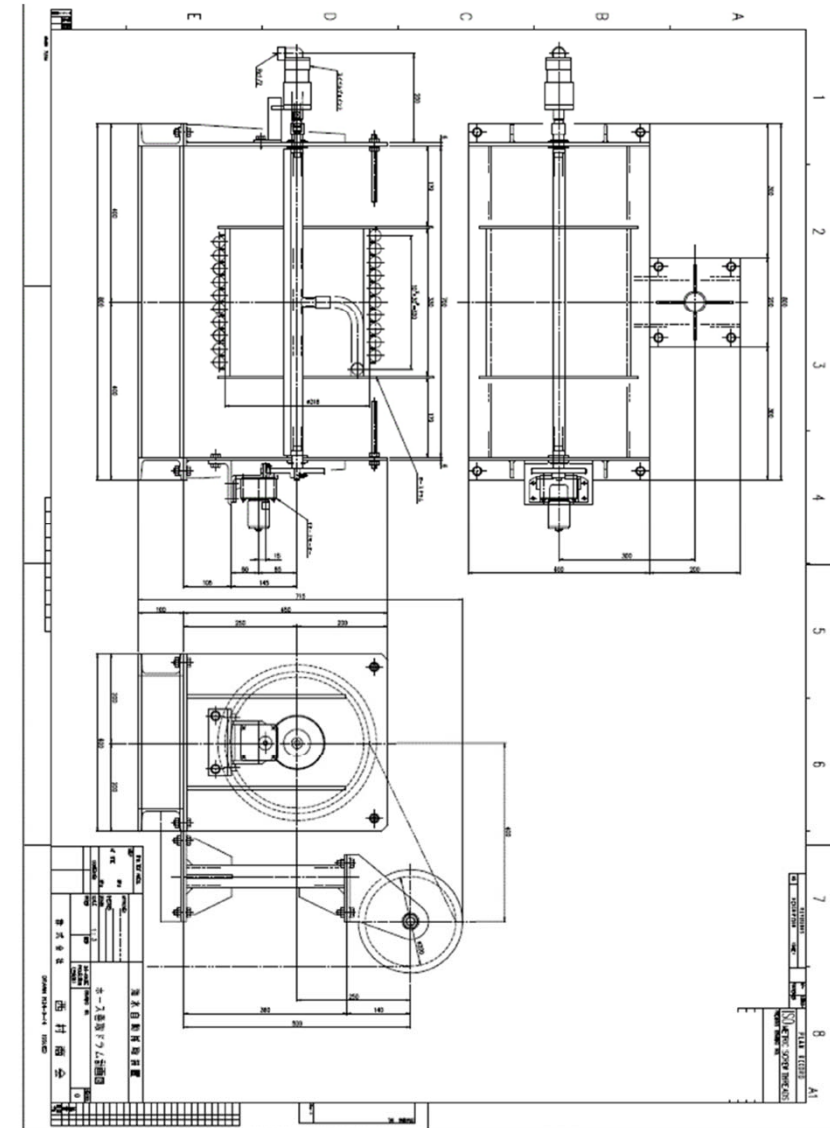


巻取ドラム仕様

駆動部

ギヤドモータ(DC12V)
減速比 1/180
回転数 30.8rpm
トルク 35kgf・cm

さらに平歯車で1/6減速
ドラムの回転数 5.13rpm
トルク 2.1kgf・m
ホースの移動速度 92mm/s





The Outline of KAUST (King Abdullah University of Science and Technology)



T +966 (2) 690-8600

Thuwal 23955-6900

Saudi Arabia





KAUST is located along the Red Sea.

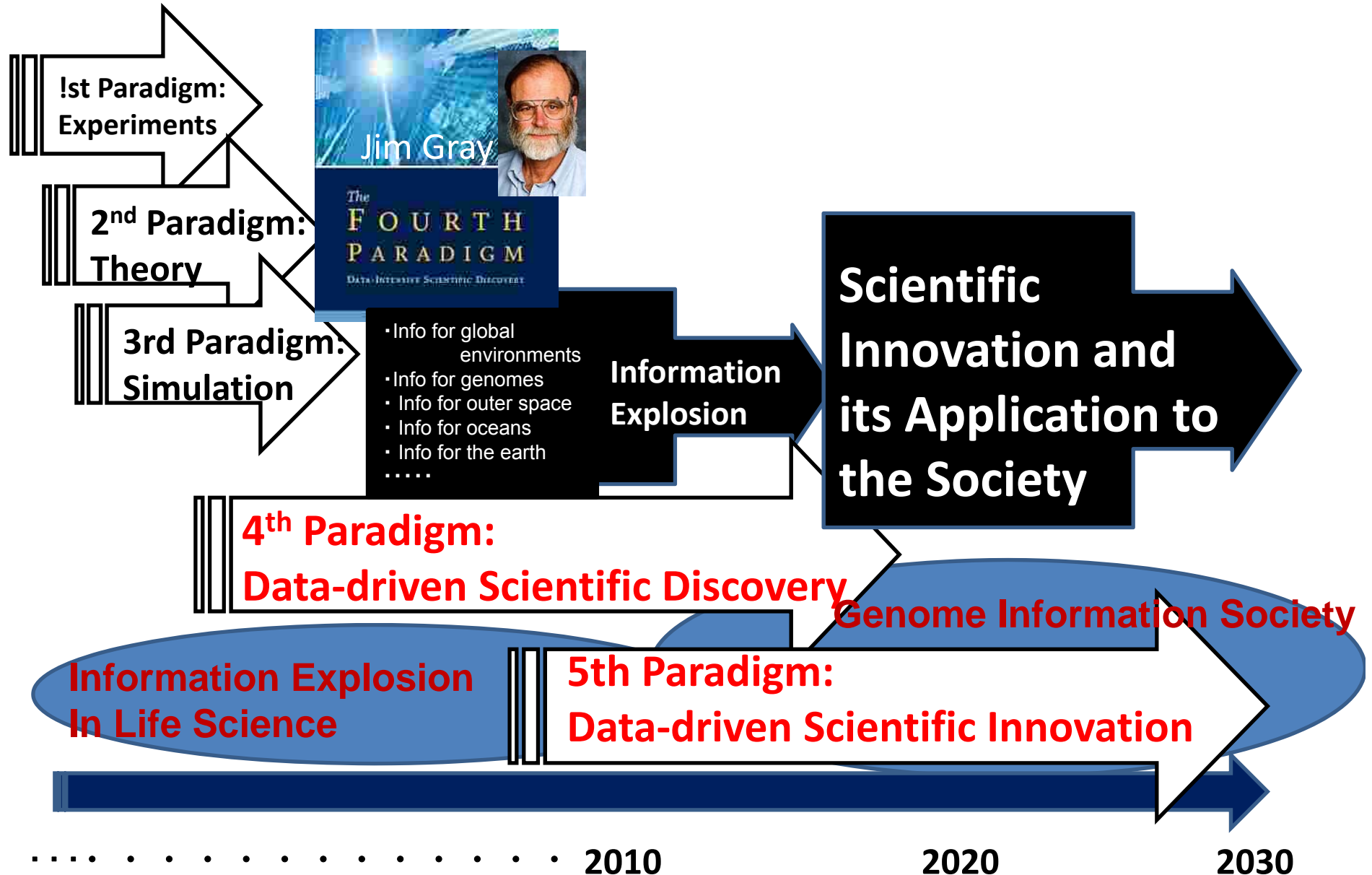


Overlooking the sea, the KAUST campus incorporates a distinctive blend of traditional architecture and modern styles.





Beyond the 4th Paradigm proposed by Jim Gray



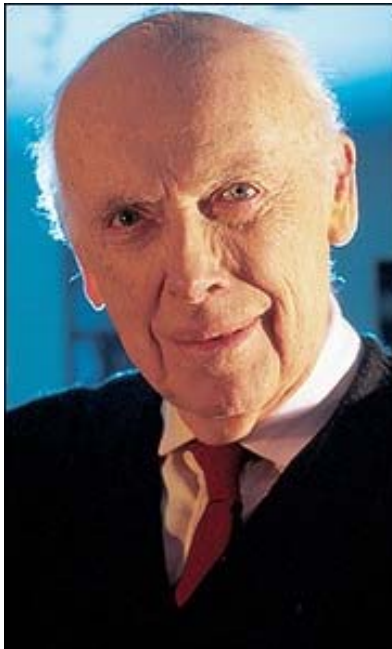
How to analyze Big Data (9)!

**“Data-driven needs working
hypothesis” is
a Key!**

特別講演

11月24日（日）14時00分～15時00分

60 years of DNA



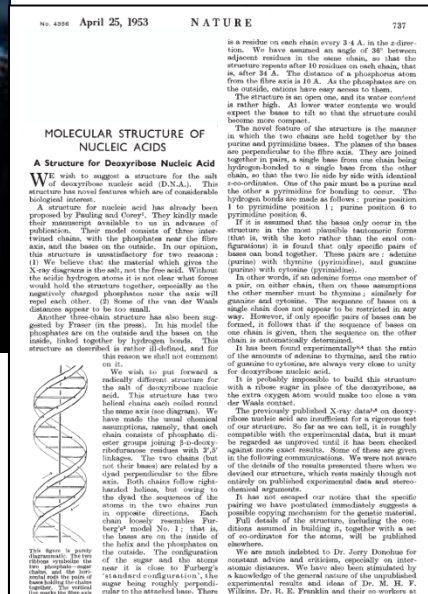
James D. Watson
(CSHL, The Nobel Prize in Physiology or Medicine 1962)

1962年度 ノーベル生理学・医学賞

略 歴

1968年～1993年 ニューヨークのコールド・スプリング・ハーバー研究所の所長
1993年～2007年 同研究所会長
1989年～1992年 NIH(国立衛生研究所)の国立ヒトゲノム研究センター初代所長

- ・ 全米科学アカデミー会員
- ・ イギリス王立協会(ロイヤルソサイエティ)会員
- ・ 大統領自由勲章
- ・ 全米科学界の栄誉とされるアメリカ国家科学賞を受ける
- ・ 海洋研究においては、ウッズホール海洋生物学研究所に在籍



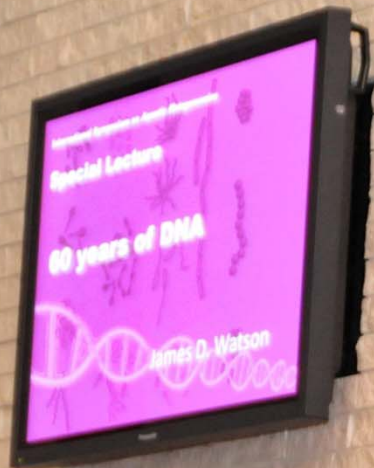




Speaker

Receiver Channels
Ch.1 日本語
Ch.2 English





Speaker





How to analyze Big Data (10)!

**“Dr. Jim Watson’s
Ten Commands” I call is
a Key!**

Collaborators

- Kazuo Ikeo (NIG)
 - Takahisa Mori (NIG)
 - Naobumi Sasaki (NIG)
 - FRA (Fishery Research Agency)
 - Ishino's group (Kyushu University)
-
- Shiho Hayakawa (British Columbia)
 - Shang Hwang (Malaysia)
 - Katsuhiko Mineta (now, KAUST)
 - Walter Gehring
(Basel, deceased on May 29, 2014)

Acknowledgements

Main work

Kazuho Ikeo (NIG)
Masafumi Nozawa (NIG)
Jung Shan Hwang (NIG)
Shiho Hayakawa (NIG)
Akiko Noda-Ogura (NIG, BIRC)
Atsushi Ogura (NIG)
Yasuharu Takaku (NIG)
Naobumi Sasaki (NIG)
Masahiro Mori (NIG)
Masakazu Ishikawa (NIG)
Masa-aki Yoshida (NIG)
Sonoko Kinjo (NIG)

Walter Gehring (Basel)
Hiroshi Suga (formerly Basel)

Planarian, Hydra, & Acidian

Kiyokazu Agata (Kyoto)
Hiroshi Shimizu (NIG)
Toshitaka Fujisawa (NIG)
Hiroaki Yamamoto (Nagaoka Bio U)

Marine Metagenomics

FRA
Kyushu University
Tokyo University
Kitasato University
Tokai University
AIST
CREST Project

Collaborators at KAUST

- Rimantas Kodzius (CBRC/BESE, KAUST)
- Hayedeh Behzad (CBRC/BESE, KAUST)
- Lujain A. Hobani (CBRC/BESE, KAUST)
- Katsuhiko Mineta (CBRC/KAUST)

- Vlad Bajic (CBRC, KAUST)
- John Archer (CBRC, KAUST)
- -----
- Shugo Watabe (Kitasato University, Japan)
- Kakehhiko Ogata (Kitasato University, Japan)

- Pei-Yuan Qian (HKUST, Hong Kong)
- Yoshizumi Ishino (Kyushu University, Japan)

DNA for World Peace !

