

Systems Biology in Supercomputing Environment

Satoru Miyano

Human Genome Center, Institute of Medical Science, The University of Tokyo BioPPN 2011 Thistle County Hotel, Newcastle upon Tyne, UK June 20, 2011



My interest on systems biology is "Cancer"

Because ...



Cause-Specific Death Rate







Die of cancer

Hack Cancer System with Supercomputer!



Caution: If you hack supercomputer, you will be arrested.





Friday, May 15 9-11PM Pacific Time





Danjuro Ichikawa – Kabuki Actor



Minselkoo HkattelameActActsee(20079)85)

NHKアーカイブス ビデオ・コンサート

天に響く歌声

アンコール 上映決定!



1985年(デビュー当時から)2004年(亡くなる前年)までの美 奈子.さんのハートフルなステージをNHKの番組からセレ クトした集大成です。アイドルからロック・クィーンそして「ミ ス・サイゴン」をきっかけにミュージカル・スター、クラシック の歌姫へと成長していく彼女の貴重な映像記録です(50 分)。

<収録曲>

~本田美奈子. 情熱のステージ~

「今夜はビートに乗れない」「1986年のマリリン」「愛の讃歌」 「新世界」「星に願いを」「アメイジング・グレイス」他

時	第1回	8月15日(金)13:00~
	第2回	8月15日(金)15:30~
	第3回	8月31日(日)13:00~
	第4回	8月31日(日)15:30~

会場 NHKアーカイブス 2F 出会いの広場(シアター)にて

それよ)このポスターの制作にみたっては、「オートレース会員資金」の補助を受けました。



Why not effective to me?

Can we hope new drugs and therapies in the future?



Cancer is Similar to Japanese Bureaucracy System





Modified: "Hanahan & Weinberg, Cell 2000.

Angiogenesis





Known mechanisms involved in Angiogenesis is very complex



Model made by using Cell Illustrator

Hallmark of Cancer



D. Hanahan and R. A. Weinberg. Cell., 100(1):57–70 Review, 2000.

The Biology of Cancer



Robert A. Weinberg

Golden Days of Molecular Biology

- Phenotypes and Genes

 –Paradigm of Molecular Biology
- Super Stars in Biology
 Novel Prize Laureates
- Dramatic History

- "The Biology of Cancer", R.A. Weinberg

〔平家物語剣之巻〕 3 軸 WA 3 1 - 5

The Tale of the Heike (13th Century, Author Unknown)

It tells the story of the rise to glory and eventual downfall of the Heike clan in the late twelfth century, a theme based on the Buddhist concept that the proud will surely be destroyed.

02

But, Very Naïve Understanding and Representation as Systems

- Draw pictures and add English narrations for biological facts
- "Systems Knowledge" which is unambiguously represented is at most causal relations among molecules
- Huge gap between "what is represented" and "what is to be represented"

Cell Illustrator Online 5.0 Java Web Start Application



Hybrid Functional Petri Net with extension (HFPNe)

Nagasaki M, Doi A, Matsuno H, Miyano S. A versatile Petri net based architecture for modeling and simulation of complex biological processes. *Genome Informatics*. 15(1):180-97, 2005.

Nagasaki, M., Doi, A., Matsuno, H., Miyano, S. Computational modeling of biological processes with Petri net based architecture. In "*Bioinformatics Technologies*" (Y.P. Chen, ed.). Springer Press. 179-243, 2005

Extension

Real (continuous)Integer (discrete)

•Object (universal) Real, Integer, Boolean, String, vector, etc.

- Dynamic cellular localization information
- Dynamic multi-cellular behavior
- Detailed molecular mechanism

User's Abilities Required for Cell Illustrator

- Biology
- Mobile phone
- Math at the level of junior high school

Place, Transition, Arc in Petri Net

	Cell Illustrator (software)					
	Original symbols of HFPNe			Examples of biological images		
Туре	Discrete	Continuous	Generic	Discrete, Continuous, and Generic		
Entity	\bigcirc	\bigcirc	\oplus	M WM		
Process						
Connector		rocess	Associatio	on Inhibitory		

Foundations of Systems Biology Using Cell Illustrator and Pathway Databases Series: Computational Biology , Vol. 13 Nagasaki, M., Saito, A., Doi, A., Matsuno, H., Miyano, S.

2009, Approx. 170 p. 145 illus. in color. With CD-ROM., Hardcover ISBN: 978-1-84882-022-7 26,95 €



2. Knowledge Representation of Dynamic Biopathways

Cell System Ontology (CSO 3.0)

Cell System Markup Language (CSML 3.0)



http://www.csml.org/

CSML 3.0 & CSO 3.0

- Native XML format for Cell Illustrator
- XML and Ontology for Biopathway
 - Modeling
 - Simulation
 - Visualization

Features of CSO 3.0

Differ from BioPAX

- <u>System-dynamics centered</u> ontology.
- The ontology is implemented with Web Ontology Language (OWL), <u>which enables</u> <u>semantic validation and provide complete and</u> <u>consistent biological pathway models.</u>

Differ from CellML and SBML

Cell System Ontology: Representation for Modeling, Visualizing, and Simulating Biological Pahtways, Euna Jeong^{*} Masao Nagasaki^{*}, Ayumu Saito, Miyano Satoru, *In Silico Biology*, 2007; 7: 0055.

Jeong E, Nagasaki M, Miyano S. Conversion from **BioPAX to CSO** for System Dynamics and Visualization of Biological Pathway. *Genome Informatics.* 2007; 18: 225-236.

Features of CSO 3.0

- CSO equips mature core vocabularies (more than 350)
- Each core vocabulary has at least one standard icon



Features of CSML 3.0

- 1. Hybrid Functional Petri net with extension (HFPNe)
- 2. Logic based descriptions, *e.g.*, temporal logic, can be defined with the format.
- 3. User can create sub-models for a model with a filtering concept.
- 4. User can define more than one views for a model, *e.g.*, gene network view, simulation view.
- 5. All terms in CSML 3.0 has the background of Ontology: Cell System Ontology (CSO) 3.0.

These parts were missing in SBML and CelIML. SMBM and CelIML are subsets of CSML3.0

3. Biopathway Layout

- Biologically sophisticated pathway layout algorithms are required
- Cell Illustrator Online 5.0 has more than 20 layout algorithms





Fast Grid Layout Algorithms for Biological Pathways



Kojima, K., Nagasaki, M.*, Miyano, S. Fast grid layout algorithm for biological networks with sweep calculation. Bioinformatics. 24(12): 1426-1432, 2008.

Kojima K, Nagasaki M, Jeong E, Kato M, Miyano S. An efficient grid layout algorithm for biological networks utilizing various biological attributes. *BMC Bioinformatics*. 2007; 8:76.

4. Transforming Pathway Databases for Cell Illustrator

XML Format

- TRANSPATH2CSML
- SBML2CSML
- CellML2CSML

Ontology Format

BioPAX2CSO

TRANSPATH to CSML

- 16 modeling rules based on Hybrid Functional Petri Net with extension (HFPNe).
- TRANSPATH (Biobase): More than 115,000 cellular events in humans, mice, and rats, collected from over 31,500 publications.
- Petri net element is incorporated with Cell System Ontology (CSO) to enable semantic interoperability of models.
- 97% of the reactions in TRANSPATH are converted into simulation-based models in CSML.

Nagasaki M, Saito A, Li C, Jeong E, Miyano S. Systematic reconstruction of TRANSPATH data into Cell System Markup Language. BMC Systems Biology. 2:53, 2008.

TRANSPATH Pathway Library Module

- More than 1,000 TRANSPATH pathways (Signal Transduction Pathway and Gene Regulatory Network) are supplied. All pathways can be loaded, edited, saved and simulated on CIO.
 - Support pathways supplied in TRANSPATH 8.4 (BIOBASE).
 - Academic user can register and use the academic version of TRANSPATH.
 - Curated 100,000 reactions and 100,000 molecules in Human and Mouse.



Pathway Parameter Search Module Data Assimilation Module → Li, Kuroyanagi et al.

• For a CIO pathway model, the module executes the user specified multiple initial conditions at once and displays the result with 2D or 3D plots.



Back to Cancer




Integrative Analysis by Top Cancer Scientists

Systems Cancer

Computational Systems Biology by Supercomputer

Credit: Lockheed Marion Corp.

Systems Cancer Research Project by MEXT (2010-2014)

Integrative Systems Understanding for Advanced Diagnosis, Therapy and Prevention



個人のがんシステムをデジタル化するには―― スーパーコンピューターインフラと最先端の生命システム解析技術を生かし、 現在のがん研究が直面している限界を超えて、がん研究の水準を飛躍的に向上・強化させる。

Computational Systems Biolog



http://systemscancer.hgc.jp



Next-Generation Supercomputer Project

1,15000,000,000JPY for 7 years (2006-2012)



SPARC64[™] VIIIfx

C Fujitsu Limited



- Design, build, and set up the Next-Generation Supercomputer with a speed of 10 PETA FLOPS.
- Over 80,000 nodes & 640,000 cores.

•Kobe(神戸)



次世代スーパーコンピュータ施設 完成イメージ図













00

Google Earth





Cancer System Hacker

Supercomputer





<u>K computer No.1</u> International Supercomputing 2011 Hamburg, June 19, 2011

80% K computer (68,544 Nodes) (8.774PFLOPS at Peak) LINPACK 8.162PFLOPS (93.0%)

Big Contributor Supercomputer System for 2009-2014



 January 2009: 75 TFLOPS at peak & 1 PB Disk Space PC Cluster (Sun Microsystems) Large Shared Memory Machine (SGI Altix) Lustre File System (Sun Microsystems)

January 2012: 225 TFLOPS at peak & 4PB Disk Space



How to Hack Cancer Systems



Integrative Systems Understanding of Cancer for Advanced Diagnosis, Therapy and Prevention

How to Hack Cancer Systems

- Digitalizing Heterogeneity/Characteristics of Cancer Systems
- Bayesian Gene Networks of Cancer Systems
- Modeling Dynamics in Cancer Systems with State Space Model
- Comparing Networks in Cancer Systems under Different Conditions
- Extracting Functional Modules in Cancer Systesm
- Software Platform eXtensible integrative Pipeline & Cell Illustrator

1. Digitalizing Cancer Systems Towards Personal Gene Networks



Teppei Shimamura, Seiya Imoto, Yukako Shimada, Yasuyuki Hosono, Atsushi Niida, Masao Nagasaki, Rui Yamaguchi, Takashi Takahashi, Satoru Miyano, PLoS ONE. 6(6): e20804, 2011.

Cancer Heterogeneity and Individual Variation

High-Throughput Technology

Microarray, ChIP-Chip, CGH array, SNP array, DNA-Seq, Exome-Seq, RNA-Seq, ChIP-Seq, ...



Patient-Specific Gene Network



Cellular System for Patient A

Cellular System for Patient B

Traditional Gene Network Estimation Problem

Knock-down Microarray Data for Sample A

	sample 1	sample 2	sample 3	••••
gene 1	1.5	5.2	1.4	••••
gene 2	5.2	6.3	0.4	••••
gene 3	3.4	9.3	0.3	••••
gene 4	2.9	0.3	6.4	••••
:		:	:	:
•	•	•	•	•

Time-series Microarray Data for Sample A

	time 1	time 2	time 3	
gene 1	1.5	5.2	1.4	
gene 2	5.2	6.3	0.4	
gene 3	3.4	9.3	0.3	
gene 4	2.9	0.3	6.4	
•	•	:	:	:
:	•	:	:	:

Gene Network for Sample A



Patient-Specific Gene Network Estimation



Corelation between 2 genes

No relation?



Change can be found if we look with a modulator







Finding Gene Relation by Data Stratification





Concept of NetworkProfiler



Modulator: Cofactor modulating relationship between genes A and B

Examples of Modulator

- Tumor progression (Stage I, Stage II, ...)
- Drug sensitivity (IC50, GI50, ...)
- Disease-free survival
- Molecular characteristics (Metastasis, EMT...)
- Pathway activity
- •

What is System?



Node: gene transcript Edge: conditional dependence (not equal to correlation)

Co-expression Network



NetworkProfiler

p (genes) × n (patients) gene expression data matrix $x_{\alpha j}$; expression value of j-th gene for α -th patient m_{α} ; modulator value for α -th patient

Structural equation model of k-th gene for α -th patient



NetworkProfiler

 $p \text{ (genes)} \times n \text{ (patients)}$ gene expression data matrix $x_{\beta j}$; expression value of *j*-th gene for β -th patient m_{β} ; modulator value for β -th patient

Structural equation model of k-th gene for β -th patient



Concept of NetworkProfiler

Purpose: estimation of **coefficients** for α -th sample $\beta_{jk\alpha}$ **Problem:** microarray for α -th sample is only **one**!! **Idea:** data integration by sample weighting



Concept of NetworkProfiler

Purpose: estimation of **coefficients** for β -th sample $\beta_{jk\beta}$ **Problem:** microarray for β -th sample is only **one**!! **Idea:** data integration by sample weighting


NetworkProfiler

Structural equation model of k-th gene for α -th patient

$$x_{\alpha k} = \sum_{j=0, j \neq k}^{p} \beta_{jk\alpha} x_{\alpha j} + \varepsilon_{\alpha k}, \quad \beta_{jk\alpha} = \beta_{jk}(m_{\alpha})$$

Elastic net-type regularized weighted loss function

$$S(\beta_{1k\alpha},\dots,\beta_{pk\alpha} \mid m_{\alpha}) = \sum_{i=1}^{n} \overline{K_{h}(m_{i}-m_{\alpha})} \left\{ x_{ik} - \sum_{j \neq k}^{p} \beta_{jk\alpha} x_{ij} \right\}^{2} + \lambda_{1} \sum_{j \neq k} w_{jk\alpha} \left| \beta_{jk\alpha} \right| + \frac{\lambda_{2}}{2} \sum_{j \neq k} \beta_{jk\alpha}^{2}$$

 $K_h(m_i - m_{\alpha})$: Gaussian kernel function (m_{α} : center, h: width) λ_1, λ_2 : regularization parameters

$$\hat{\boldsymbol{\beta}}_{k\alpha} = \left(\hat{\beta}_{1k\alpha}, \dots, \hat{\beta}_{pk\alpha}\right)^T = \underset{\boldsymbol{\beta}_{k\alpha}}{\operatorname{arg\,min}} S(\beta_{1k\alpha}, \dots, \beta_{pk\alpha} \mid m_{\alpha})$$

The neighborhood samples in terms of modulator have similar network structure

Model Selection

Performance of NetworkProfiler \rightarrow Selection of λ_1, λ_2 and **h**

Select λ₁ and λ₂ based on WAICc (Shimamura et al., 2010b)
Select h based on cross-validation

$$S^{(-i)}(\beta_{1k\alpha},...,\beta_{1k\alpha} \mid m_{\alpha}) = \sum_{i \neq \alpha} K_{h}(m_{i} - m_{\alpha}) \left\{ x_{ik} - \sum_{j \neq k}^{p} \beta_{jk\alpha} x_{ij} \right\}^{2} + \lambda_{1} \sum_{j \neq k} w_{jk\alpha} \left| \beta_{jk\alpha} \right| + \frac{\lambda_{2}}{2} \sum_{j \neq k} \beta_{jk\alpha}^{2}$$
$$\hat{\boldsymbol{\beta}}_{k\alpha}^{(-i)} = \arg \min \left\{ S^{(-i)}(\beta_{1k\alpha},...,\beta_{1k\alpha} \mid m_{\alpha}) \right\}$$
$$CV_{k}(h) = \sum_{\alpha=1}^{n} \left\{ x_{\alpha k} - \sum_{j \neq k} \hat{\beta}_{jk\alpha}^{(-\alpha)} x_{\alpha j} \right\}^{2}$$

Select *h* minimizing $CV_k(h)$

Epithelial-Mesenchymal Transition (EMT)

- Key developmental remodeling program, where cells alternate between Epithelial-like and Mesenchymal-like phenotypes
- Relate to tumor grade and metastasis
- Contribute to increasing in drug resistance Chua

Chua et al. (2008)



System related to EMT is a "black box"

Modulator for EMT

We selected coherent 50 genes from 121 EMT signature genes to define the modulator for EMT (EEM, Niida et al., 2009)



Modulator for EMT

We selected coherent 50 genes from 122 EMT signature genes to define the modulator for EMT (EEM, Niida et al., 2009)

Signature-based hidden modulator extraction algorithm

1. Selected 122 genes labeled "EMT UP", "EMT DN", "JECHLINGER EMT UP", and "JECHLINGER EMT DN"from Molecular Signatures Database v2.5 ([6];

<u>http://</u>

www.broadinstitute.org/ gsea/msigdb/index.jsp).

- 2. Then, narrowed the set to 50 functionally coherent genes with $p < 10^{-5}$ by using the extraction of expression module (EEM).
- Computed the first principal component of these 50 genes as hidden values of the EMT-related modulator



Elucidating Systems Responsible for EMT

Input

- Transcriptome data of 762 cancer cells (22,283 transcripts=22277 mRNA+581 miRNA)
- EMT Modulator





miRNAs highly expressed in Epithelial



System Changes Related to EMT

- Functional loss of E-cadherin = a hallmark of EMT
- Focus on regulators of E-cadherin

Epitherial-Like Cell

Mesenchymal-Like Cell



Upstream Regulatory Changes of E-cadherin

Coefficient profiles of E-cadherin regulators through EMT



regulator	$_{\mathrm{type}}$	regulatory effect change	Evidence
IRF6	А	101.04	
miR-141	А	87.58	[8]
GRHL2	Α	64.13	
ZEB1 (SIP1)	Ι	50.72	[9]
LSR	Ι	46.89	
miR-200b	А	31.55	[8]
KLF4	А	26.28	[10]
OVOL2	А	22.08	
miR-200a	А	17.70	[8]
FOXA2	А	17.26	[11]
TCF4 (E2.2)	Ι	14.15	[12]
ELF3	А	13.58	
ZNF217	Α	13.53	
MYB	А	12.50	
KLF5	А	12.42	
miR-192	А	12.30	[13, 14]
FOXA1	А	11.69	[11]
ZNF165	А	11.39	
NKX2-1	А	11.21	
HNF1B	А	11.08	
TFE3	А	11.01	
ZEB2 (δ EF)	Ι	10.66	[15]
TRIM29	Ι	9.87	
SNAI2	Ι	9.74	[16]

Table 1. 25 top-ranked regulators of E-cadherin for the change in the regulatory effect change among the EMT with published evidence

A: Activator I: Inhibitor

Upstream Regulatory Changes of E-cadherin



(c). The green and red colors indicate epithelial- and mesenchymal-like cells, respectively.







miR-141, ZEB1, and E-cadherin

- NetworkProfiler revealed regulatory changes in *miR-141*, *ZEB1*, and E-cadherin.
- Specifically, it suggested that decreased expression of *miR-141* in mesenchymal cells disrupts the negative feedback loop between*miR-141* and *ZEB1*, which would allow *ZEB1* to decrease the expression of E-cadherin during the EMT.

- We predicted 45 EMT-dependent putative master regulators that control sets of genes involved in cell adhesion, migration, invasion and metastasis, namely,
- 17 of which are in the downstream of TGFB1, a master switch of the EMT, in our prediction.

regulator	$_{\mathrm{type}}$	regulatory effect change	Evidence
IRF6	А	101.04	
miR-141	А	87.58	[8]
GRHL2	Α	64.13	
ZEB1 (SIP1)	Ι	50.72	[9]
LSR	Ι	46.89	
miR-200b	А	31.55	[8]
KLF4	А	26.28	[10]
OVOL2	А	22.08	
miR-200a	А	17.70	[8]
FOXA2	А	17.26	[11]
TCF4 (E2.2)	Ι	14.15	[12]
ELF3	Α	13.58	
ZNF217	А	13.53	
MYB	А	12.50	
KLF5	А	12.42	
miR-192	А	12.30	[13, 14]
FOXA1	А	11.69	[11]
ZNF165	А	11.39	
NKX2-1	А	11.21	
HNF1B	А	11.08	
TFE3	А	11.01	
ZEB2 (δ EF)	Ι	10.66	[15]
TRIM29	Ι	9.87	
SNAI2	Ι	9.74	[16]

Table 1. 25 top-ranked regulators of E-cadherin for the change in the regulatory effect change among the EMT with published evidence

A: Activator I: Inhibitor

A novel regulator KLF5 of EMT

- Krueppel-like factor 5 (KLF5) from a list of the remaining candidate regulators and conducted *in vitro* validation experiments.
- As a result, we found that knockdown of KLF5 by siRNA significantly decreased Ecadherin expression and induced morphological changes characteristic of EMT.



Relapse risk



Patient Samples (Microarray Data)



Relapse risk score





Network of the lowest risk patient

Network of the highest risk patient

Differences of Hubness Suggest Key Genes



System-Oriented Personalized Medicine

Networks

EMT/Prognosis/Metastasis/etc.





Thorough Analysis Requires 500 Times

25 yeas with 1000 cores 10 days with 1,000,000 cores Prognosis Metastasis Resistance EMT

Prologue

Mapmakers in Systems Biology

伊能忠敬 Tadataka Inoh A man who walked 40,000,000 steps



 伊能忠敬は、江戸幕府の事業として、1800年 から1816年にかけて全国を歩いて測量をし、 1821年に「大日本沿海輿地全図」が幕府に納 められたといいます、伊能忠敬はその完成を 見ずに1818年に死去ましたが、その後、仕上 げの編纂作業が行われ、全部で21年の歳月 をかけてこの地図は完成しています。

2001年のNHKの正月時代劇「四千万歩の男・伊能忠敬」 (原作:井上ひさし「四千万歩の男」(講談社))として放映



The Mapmakers

The story of the great pioneers in cartography – from antiquity to the space age

•John Noble Wilford



2. Mining Gene Networks from Gene Expression Profiles Mapmaking of Molecular Networks

- Bayesian Network Gene Networks
 - Gene knock-down/knock-out
 - Various shocks
 - Time-course data

Yoshinori Tamada;Seiya Imoto;Masao Nagasaki;Satoru Miyano

2. Bayesian Networks + Nonparametric Regression

- Gene network: model for transcriptome level gene-gene regulation using directed graphs.
- We want to estimate gene networks from high-throughput biological data e.g. gene expression data.

-1.54

-2.1

1.23

1.44



What we wanted to do



expression data

Bayesian Network and Nonparametric Regression

Network of 521 genes constructed from 120 yeast microarrays obtained by disrupting 120 genes, where 78 of them are transcription factors.

🔽 Gene name		
🗖 Systematic name		
🗖 Gene comment		
Edge comment		
Search from		
I items		
C selected items		

G.NET Application version 0.1

- Imoto, S., Goto, T., Miyano, S. Estimation of genetic networks and functional structures between genes by using Bayesian network and nonparametric regression. Pacific Symposium on Biocomputing. 7:175-186, 2002.
- 2. Imoto, Kim, Goto, Aburatani, Tashiro, Kuhara, Miyano S. Bayesian network and nonparametric heteroscedastic regression for nonlinear modeling of genetic networkJ. *Bioinformatics and Comp. Biol.*, 1(2), 231-252, 2003



Nonlinear Bayesian network model

$$f(x_{i1},...,x_{ip};\boldsymbol{\theta}_{G}) = \prod_{j=1}^{p} f_{j}(x_{ij} | \mathbf{p}_{ij};\boldsymbol{\theta}_{j}),$$

$$f_{j}(x_{ij} | \mathbf{p}_{ij};\boldsymbol{\theta}_{j}) = \frac{1}{\sqrt{2\pi\sigma_{j}^{2}}} \exp\left\{-\frac{(x_{ij} - \mu_{ij})^{2}}{2\sigma_{j}^{2}}\right\}$$

$$\mu_{ij} = m_{1}(p_{i1}^{(j)}) + \dots + m_{q_{j}}(p_{iq_{j}}^{(j)})$$

$$= \sum_{k=1}^{q_{j}} \sum_{m=1}^{M_{jk}} \gamma_{mk} b_{mk}^{(j)}(p_{ik}^{(j)})$$

Criterion for Selecting Good Networks BNRC Score

Bayesian Network and Nonparametric Regression Criterion

$$\begin{aligned} \text{BNRC}(G) &= -2\log \pi_G \int \prod_{i=1}^n f(\mathbf{x}_i; \mathbf{\theta}_G) \pi(\mathbf{\theta}_G \mid \boldsymbol{\lambda}) d\mathbf{\theta}_G \\ &= -2\log \pi_G - r\log(2\pi n^{-1}) \\ &+ \log \left| J_{\boldsymbol{\lambda}}(\hat{\mathbf{\theta}}_G) \right| - 2nl_{\boldsymbol{\lambda}}(\hat{\mathbf{\theta}}_G \mid \mathbf{X}_n) \end{aligned}$$

We choose the graph that minimizes the value of the BNRC score.



Optimal to Locally Optimal WR: Optimal Bayesian Networks of 32 Nodes April 2010



Parallel Computing with 8192 Cores



Genome-Wide Bayesian Network Computation



Anti-cancer Drug Response Gene Network of Melanoma

Dynamic Bayesian + Nonparametric Regression

*

▶
 ★

~

 \mathbf{V}

Melanoma A-375 + Taxol (Paclitaxel)
 Inferring and chasing the changes of gene networks of 2,000 genes for 24 hrs at 14 time points (triplicate at each time)
 1 hour computation using 1024 cores

 Mouse position: 895 : 5
 開く Done.
 Selection
 1865:0:0:3824
 15:44:14
 値 5244M のうちの 1900M

 t =
 1
 2
 3
 4
 5
 6
 7
 8


















Gene Networks of Small Airway Epithelial Cell and Gefitinib State Space Model -

Growth Factor Signaling Systems Identify Critical Genes for Survival Prediction in Early Stage Lung Adenocarcinoma

Yamaguchi, R., Imoto, S., Yamauchi, M., Nagasaki, M., Yoshida, R., Shimamura, T., Hatanaka, Y., Ueno, K., Higuchi, T., Gotoh, N., Miyano, S.



State Space Model: $x_n = Fx_{n-1} + v_n \in \mathbb{R}^k$ System Model $y_n = Hx_n + w_n \in \mathbb{R}^p$ Observation Model

High-Dimensional Short Time Series Data:

$$Y = \left\{ y_1, \cdots, y_{N_{\text{obs}}} \right\} \qquad N_{\text{obs}} << p \approx 10^3$$

System Estimation with Dimension Reduction:

 $\dim(x_n) = k < \dim(y_n) = p$

$$L(\theta) = \int p(Y, X | \theta) dX$$
$$H^{T} R^{-1} H = \Lambda = \operatorname{diag}(\lambda_{1}, \dots, \lambda_{k})$$

Gene Expression Prediction:

$$y_{n|n-1} = H \int x_n p(x_n \mid y_1, \cdots, y_{n-1}, \theta) dx_n$$

Module-Based Gene Network Estimation:

 $R^{-1/2}(y_n - w_n) = \Psi R^{-1/2}(y_{n-1} - w_{n-1}) + R^{-1/2}Hv_n^{1/2}$

Epidermal Growth Factor Receptor Pathway



State Space Model for Inferring Transcriptional Module Networks from Time-Course Gene Expression Data



SSM Application Predicting Differences in Gene Regulatory Systems

- Focus: EGFR pathway
- Data
 - Time Course Gene Expression Microarray Data (20K)
 - 19 time points during 48 hours
 - Human Small Airway Epithelial Cell (SAEC)
 - Two Conditions (Different Drugs)
 - EGF Stimulation (Control Data)
 - EGF + Gefitinib Dosed (Case Data)



Selection 1500 Genes for Network Analysis

а



Literature based knowledge

Time-Course of Drug-Stimulated Human Normal Lung



Gene Selection for SSM Analysis -----> 1500 Genes: Literature DB (Ingenuity) + Variation Filter (variance)

Selection of System Dimension: k Parameter Estimation: θ

→ Learning SSM with EGF data (Control's System)

Learning SSM with EGF data (Control's k = 8 minimizes prediction errors for

hold-out sample.

Meta Analysis of Transcriptional Modules

> Construction of Gene Network

Prediction of EGF-Gefitinib data by EGF-Learned SSM Obtaining *p*-value for each Gene



Genes with small p-values are considered to be induced differential regulations by Gefitinib

Differentially-Expressed and Differentially-Regulated Genes



Differential Regulations by Drug Dosing



Predicting Case Data by Control's System to Discriminate the Two Situations

A Thought Experiment: If we know the Control's System, we use it to predict the Case data.



We use a statistical model for inferring gene regulatory systems.

Strategy to Predict Differentially Regulated Genes

1. Train SSM by CTRL time-course data and Estimate an CTRL-SSM System



2. Predict CASE time-course data by the CTRL-SSM System



Prediction of EGF Data and EGF-GFT Data by SSM of EGF System

EGF Data Prediction: Check for Model Accuracy EGF+GFT Data Prediction: Exploration for Differences between Systems

Good Prediction

Bad Prediction



- X : EGF obs (Control)
- O: EGF+GFT obs (Case)
- ---- : EGF pred by SSM(EGF)
 - : EGF+GFT pred by SSM(EGF)

Differentially Regulated Genes

Examples: Selected Diff Reg Genes

Similarly Regulated Genes



Estimated Gene Networks with SSMs





ARTICLES



Gene expression-based survival prediction in lung adenocarcinoma: a multi-site, blinded validation study

Director's Challenge Consortium for the Molecular Classification of Lung Adenocarcinoma:¹¹ Kerby Shedden^{2,3,17}, Jeremy M G Taylor^{3,4,17}, Steven A Enkemann^{5,17}, Ming-Sound Tsao^{6,17}, Timothy J Yeatman^{5,17}, William L Gerald^{7,17}, Steven Eschrich^{5,17}, Igor Jurisica^{6,17}, Thomas J Giordano⁸, David E Misek³⁹, Andrew C Chang^{3,9}, Chang Qi Zhu⁶, Daniel Strumpf⁶, Samir Hanash³, Frances A Shepherd⁶, Keyue Ding¹⁰, Lesley Seymour¹⁰, Katsuhiko Naoki¹¹, Nathan Pennell¹¹, Barbara Weir¹¹, Roel Verhaak¹¹, Christine Ladd-Acosta¹², Todd Golub¹², Michael Gruidl⁵, Anupama Sharma⁵, Janos Szoke⁷, Maureen Zakowski⁷, Valerie Rusch⁷, Mark Kris⁷, Agnes Viale⁷, Noriko Motoi⁷, William Travis⁷, Barbara Conley¹³, Venkatraman E Seshan^{14,17}, Matthew Meyerson^{11,12,17}, Rok Kuids^{3,17}, Kevin K Dobbin^{15,17}, Tracy Lively^{16,17}, James W Jacobson^{16,17} & David G Beer^{3,9,17}

Although prognostic gene expression signatures for survival in early-stage lung cancer have been proposed, for clinical application, it is critical to establish their performance across different subject populations and in different laboratories. Here we report a large, training-testing, multi-site, blinded validation study to characterize the performance of several prognostic models based on gene expression for 442 lung adenocarcinomas. The hypotheses proposed examined whether microarray measurements of gene expression either alone or combined with basic clinical covariates (stage, age, sex) could be used to predict overall survival in lung cancer subjects. Several models examined produced risk scores that substantially correlated with actual subject outcome. Most methods performed better with clinical data, supporting the combined use of clinical and molecular information when building prognostic models for early-stage lung cancer. This study also provides the largest available set of microarray data with extensive pathological and clinical annotation for lung adenocarcinomas.

In the United States and in many Western countries, lung cancer represents the leading cause of cancer-related death¹. The 5-year, overall survival rate is 15% and has not improved over many decades. This is mainly because approximately two-thirds of lung cancers are discovered at advanced stages, for which cure by surgical resection is no longer an option. Furthermore, even among early-stage patients who are treated primarily by surgery with curative intent, 30-55% will develop and die of metastatic recurrence, Recent multinational dirical trials (IALT [BR10, ANITA, UFT, IACE) conducted in several continents have demonstrated that adjuvant chemotherapy significantly improves the survival of patients with early-stage (IB-II) disease². Nevertheless, it is dear that a proportion of patients with stage I disease have poorer prognosis and may benefit significantly from adjuvant chemotherapy. whereas some with stage II disease with relatively good prognoses may not benefit significantly from adjuvant chemotherapies. It remains possible, however, that the latter patients could derive additional benefit from adjuvant targeted therapies2-4. Therefore, there is an urgent need to establish new diagnostic paradigms and validate in clinical trials methods for improving the selection of stage I-II patients who are most likely to benefit from adjuvant chemotherapy.

Global gene-expression profiling using microarray technologies has helped to improve our understanding of the histological heterogeneity of non-small cell lung cancer (NSCLC) and has identified potential biomarkers and gene signatures for dassifying patients with significantly different survival outcomes5-11. However, the performance and general applicability of published dassifiers has not been easy to establish because of small numbers of subjects examined and inclusion of heterogeneous tumor types. Furthermore, there have not been uniform criteria for sample inclusion, annotation, sample processing and data analyses. To address these concerns and to generate a large microarray database of NSCLC samples that have been collected and studied using a common protocol¹², we conducted a large retrospective, multi-site, blinded study. The study included a blinded validation step to characterize the performance of several newly developed prognostic models using a total of 442 lung adenocarcinomas, the specific type of lung cancer that is increasing in incidence¹³.

To ensure scientific validity of the results, subject samples along with all relevant clinical, pathological and outcome data were collected by investigators at four institutions using data from six lung-cancer treatment sites with subject indusion criteria defined a priori. Gene

The complete list of affiliations appears at the end of the article. Correspondence should be addressed to J.W.J. (jacobsonj@ctep.ncl.nih.gov) or D.G.B. (dgbeer@umich.edu).

Received 8 January; accepted 12 June; published online 20 July 2008; doi:10.1038/nm.1790

Survival Predictions for Lung Cancer Patients with Gene Sets Identified by SSM analysis for SAEC data



Classifier: Risk Score Function

- Partial Cox Regression Model
 - Hazard Function: $\lambda(t, X_i) = \lambda_0(t) \exp[f(X_i)]$
 - Risk Score Function $f(X_i^{\text{Valid}}) = \sum_{i=1}^{p} \beta_j^{*\text{Train}} \left(X_{ij}^{\text{Valid}} - \overline{x}_j^{\text{Train}} \right)$ where $X_i^{\text{Valid}} = \begin{bmatrix} X_{i1}^{\text{Valid}}, L, X_{ip}^{\text{Valid}} \end{bmatrix}^T$ The p gene expressions for the *i*th patient in a validation set

The *i*th patient is classified into a high risk group when or a low risk group when

$$f(X_i^{\text{Valid}}) > 0$$
$$f(X_i^{\text{Valid}}) \le 0$$

Li and Gui, 2004

Survival Predictions for Lung Cancer Patients with Gene Sets Identified by SSM analysis for SAEC data



4. Gene Networks of Lung Cancer

Gene Networks of Gefitinib Sensitive PC9 and Gefitinib Resistant PC9GR2

PC9GR2 (Lung Cancer Cell Line with Gefitinib resistance) PC9GRM2 (hours) ne hour 24 27 30 33 36 39 43 48 Andre Fujita: Statistician **RNA** Sampling Yoshinori Tamada: Computer Scientist 24 hr EGF: 0 hr ~ 24 hr

EGF+Gefitinib : 0 hr ~ 24 hr

Gefitinib : 0 hr ~ 24 hr

mRNA

Total: 102 time points

(26x4 - 2 = 102)






















































Dynamic Bayesian Network + Nonparametric Regression







5. Building Data Analysis and Simulation Pipeline at ONE STOP XiP (eXtensible integrative Pipeline)

Systems Biology integrative Pipeline



NetComparator

Shimamura, T., Imoto, S., Yamaguchi, R., Nagasaki, M., Miyano, S. Inferring dynamic gene networks under varying conditions for transcriptomic network comparison. Bioinformatics. 26(8):1064-1072, 2010.



Regularized Weighted Recursive Elastic Net





Plot



Difference between Gefitinib Sensitive Lung Cancer and Resistant Lung Cancer



E2F1 is knows as a TF regulating apoptosis and cell cycle

Gefitinib Sensitive Lung Cancer

Gefitinib Resistant Lung Cancer

NCOA3: a nuclear receptor coactivator that interacts with nuclear hormone receptors toenhance their transcriptional activator functions.

A New Paradigm for Understanding

concer

© 2009 Here and Now Books



Super-early biomarker

Prediction of efficacy and relapse

New molecular targets

Mechanism of drug resistance

Systems Understanding of Cancer