

オミックス医療の現状と今後の課題

東京医科歯科大学 データ科学推進室

東北大学 東北メディカル・メガバンク機構

田中 博



ゲノム・オミックス医療の現状

バイオテクノロジーの
急速な進展による

「ビッグデータ医療」時代の到来

オミックス医療学への多大な影響 ビッグデータ時代の到来

- (1) 次世代シーケンサ (Clinical Sequencing) を始めとする「ゲノム/オミックス医療」における網羅的分子情報収集/蓄積
- (2) **Biobank/ゲノムコホート普及**による分子・環境情報の蓄積
- (3) **モバイルヘルス(mHealth)** によるWearable センサの連続計測による生理データの蓄積 (unobstructed monitoring)

コストレスで良質なデータが大量に収集可能

治療医学の**的確性の飛躍的進展** 「精密医療」
医療の国民レベル・生涯ヘルスケアの**進展**

ゲノム・オミックス医療の2つの流れ

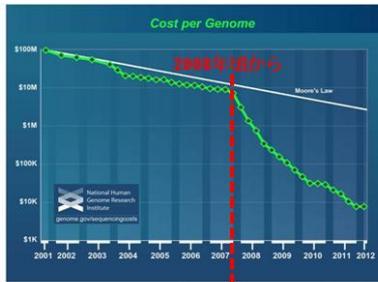
米国での流れ

- 次世代シーケンサの急激な発展による「シーケンス革命」からの怒濤の展開（2010から）
- 「治療医学」レベル質的向上のためにゲノム情報を取り入れた臨床実装の推進 3つのゲノム・オミックス医療
 - 稀少疾患の原因遺伝子変異の同定
 - がんのドライバー遺伝子変異の同定と分子標的薬の選択
 - 薬剤代謝酵素の多型性の同定と個別化投与

欧州での流れ

- 社会福祉国家の理念より国民医療（医療の国民レベル）の向上
- 「予防医学」レベル質的向上のためにゲノム情報を取り入れたバイオバンク推進
- 大規模前向きpopulation型バイオバンク/ゲノム・コホートの確立
 - 遺伝的素因だけでなく環境要因（生活習慣）との相互作用を解明し、「ありふれた疾患」発症を予測し、これに基づいて個別化予防する。
 - 疾患を発症前に対応して発症を防ぐ「先制医療(preemptive medicine)」や「予測医療 (predictive medicine)」の実現を目的

米国ゲノム・オミックス医療の流れ



DNA Sequencing Cost: the National Human Genome Research Institute

シーケンス革命 2007/8

2005~ NGS 454 (LS,Roche)
2007/8~454, Solexa (Illumina), SOLiD (LT,TF)
シーケンス革命



	HiSeq2500	Ion Proton
本体価格	約1億円	約3500万円
モード / チップ	ハイアウトプット	ラビッドラン
解析時間	11日	27時間
リード長 (bp)	2 x 100	2 x 150
データ産出量 (Gb)	約600	約120
試薬コスト (ヒト1人全ゲノム)	数十万円	不可 エクソームのみ

急速な高速化と廉価化
ヒトゲノム解読計画13年,3500億円
⇒1日,10万円



オバマ大前統領 Precision Medicine Initiativeを開始、2015年1月 大統領一般年頭教書演説

先陣争いの時代

第一期

ゲノム多型性の認識
Hapmap 計画(2002)
GWAS研究など

薬剤代謝酵素の多型性の判別・電子カルテで警告・Preemptive PGx
Vanderbilt大病院

2005~ NGS 登場
(454, Solexa, SOLiD)
2007/8~
シーケンス革命

Undiagnosed Genetic Diseaseの原因遺伝子POC同定
MCW小児病院

国際がんコンソーシアム開始
ICCG (2008年)
2011頃からがん変異成果報告

Cancer Driver Geneの同定と抗がん剤治療
Dana Faber CC

ゲノム・オミックス医療の臨床実装の普及
ゲノム・オミックス情報のビッグデータの出現

第二期

国家政策の時代

ゲノム医療の国家的取組み
NIH "BD2K" 計画・各種ゲノムコンソーシアム開始

オバマ大統領 年頭教書
Precision Medicine initiative 政策の発表

第三期

精密医療普及期

100万人コホート:バンダービルト大学設開始(2016-2020)
NCI "National Cancer MoonShot" 10年計画開始
各州でのプレジジョン医療計画開始(カリフォルニア、ペンシルバニア)

2007年

2009年

2010年

2011年

2012年

2013年

2014年

2015年

2016年

2017年

個別化医療から Precision Medicine

個人の遺伝素因・環境要因に合わせた (tailored) 医療
One size fits for all の Population 医療とは異なる

趣旨：基本は、個別化医療 Personalized Medicine の概念と変わらないが
目的は診断/治療の個人化ではなく層別化を明確化

概念の拡張：Personalized Medicineが標榜された時から10数年経っている

医療ビッグデータ時代の到来による個別化医療の拡張

(1) 遺伝素因 X 環境(生活習慣)要因のスキーマ重視

遺伝要因(SNPや変異：Genome)だけでなく環境・生活習慣要因(Exposome)の重視
疾患発症は2つの要因の相互作用と明快に強調。電子カルテの臨床表現型
(Clinical Phenome)情報の重要性認識。

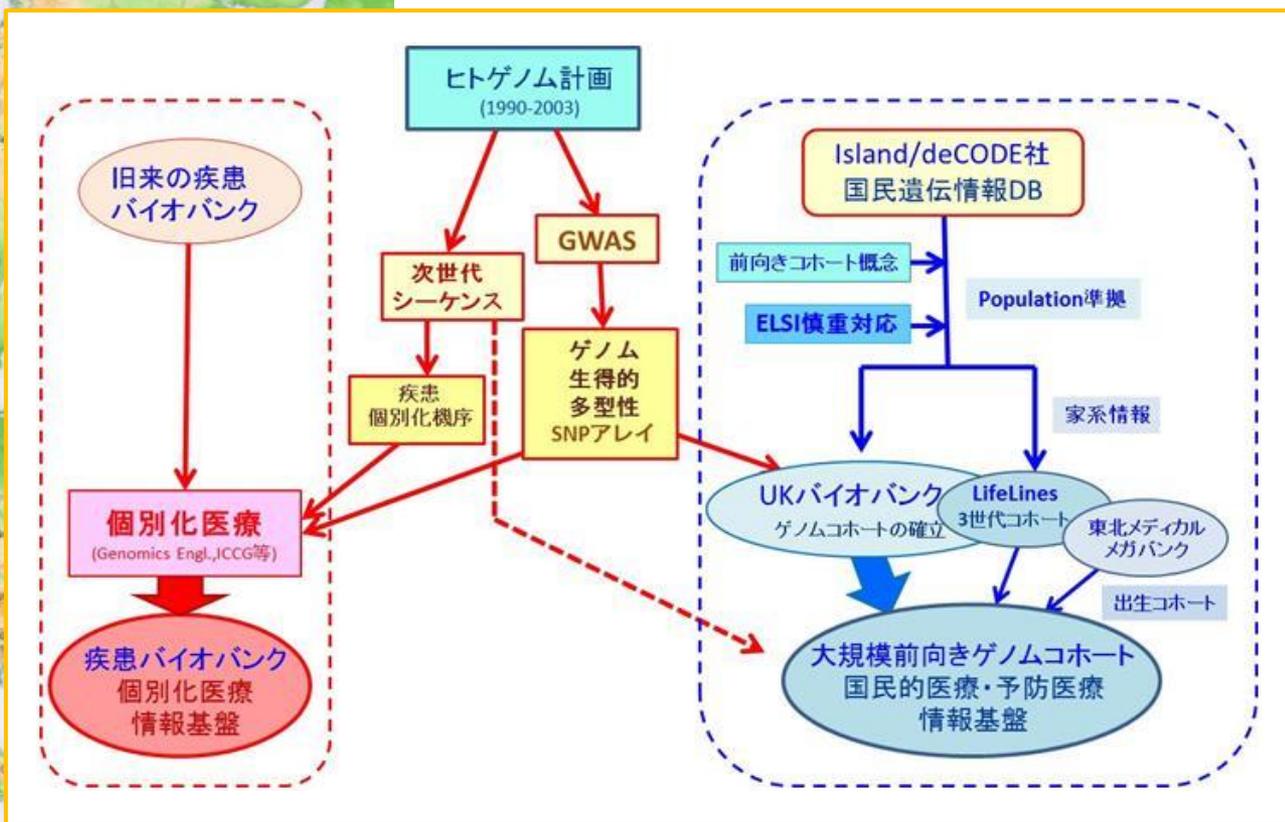
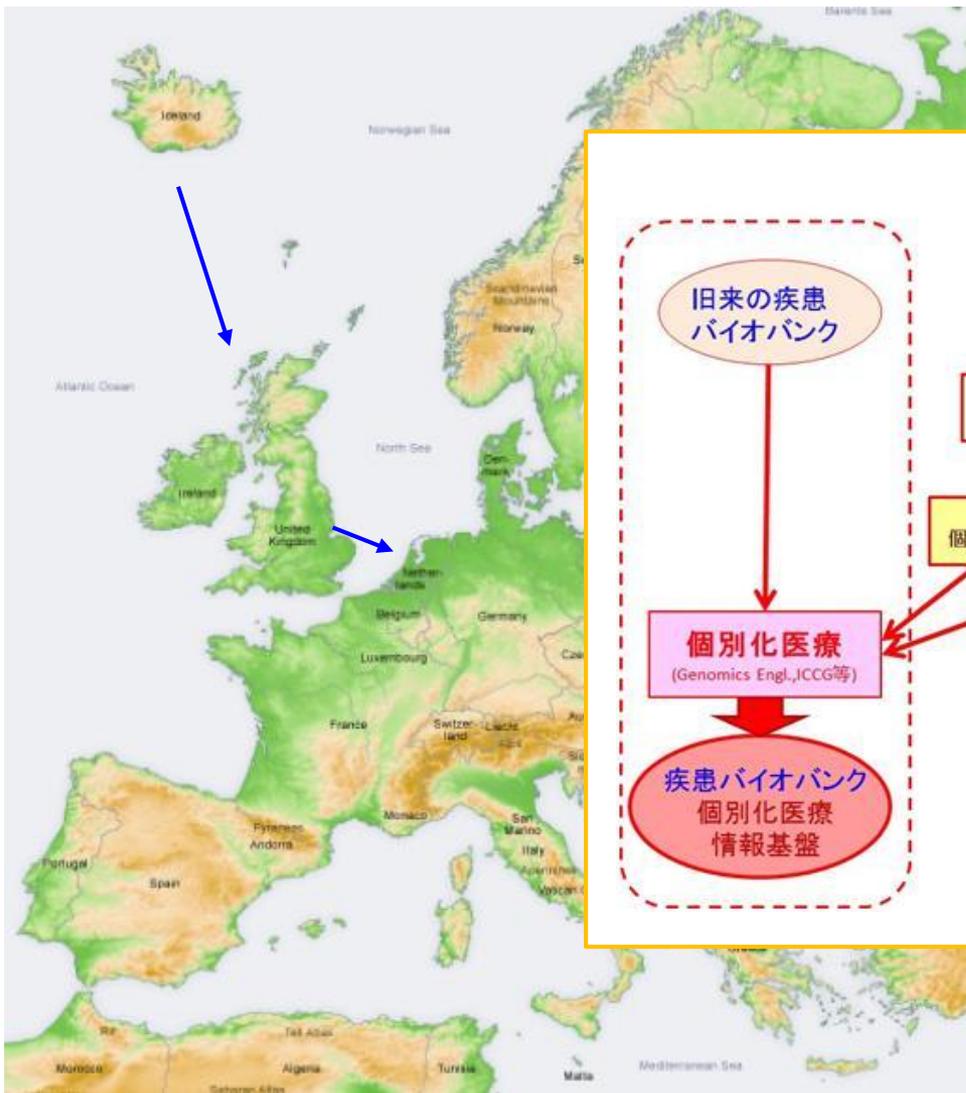
(2) ゲノムコホート・Biobankの重視

Precision Medicineを実現する「情報基盤」として、ゲノムコホート/Biobankが
必要であることを認識。Real world dataの重視

(3) 日常生理モニタリング情報の包摂

モバイルヘルス(mHealth)・wearableセンサーによる大量継続情報収集の重視

第2の流れ 欧州のバイオバンクの普及



ビッグデータ医学/医療の第2の流れ

Biobankとゲノムコホートの世界的興隆

バイオバンクの目的・機能の変化

- 従来は**稀少疾患組織標本**や臨床研究の**資料保存**、近年は**ゲノム医療の基盤**としての役割が認識され、**世界的に普及**
- **疾患BioBank** : **ゲノム・オミックス個別化医療/創薬**の情報基盤 :
 - 疾患罹患患者の網羅的分子情報（ゲノムなど）とそれに対応する**臨床表現型情報**の収集。
 - 疾病の分子機序や治療戦略、**予後予測**、**創薬科学**への貢献
 - 従来の疾患バイオバンクが**個別化医療**の概念により**変革**、**個別化医療の情報基盤**としての役割
- **Population型BioBank** : **国民医療レベルの向上**、**予防**の情報基盤 :
 - 「健常者」前向きコホート。調査開始時の網羅的分子情報（ゲノム）と**臨床・環境情報**（exposome）を集めて、**長期間（生涯）を追跡するゲノム・コホート**
 - 主に**遺伝子素因情報**も含めた「**ありふれた病気**」の**疾患の発症リスク予測**、**重症化予測**

欧米のBiobank

- **英国 UK biobank**
 - 50万人の健常者。40~69歳（2006-2010, 62Mポンド）、追加調査（2011-16, 25Mポンド）
 - 健診データ（血液・尿・唾液サンプル、生活情報）とゲノム情報（SNPアレイを集め、健康医療状況を追跡する。その淵源は、アイスランド、deCODE社の「国民遺伝子情報データベース」プロジェクト
- **英国 Genomics England**,
 - 2013開始、2017年までに10万人のゲノム配列収集。全ゲノム次世代シーケンス
 - 最初の対象は稀少疾患（患者・家族）、がん患者、最初はEnglandのみ。企業とのコンソーシアム
- **欧州 BBMRI** (Biobanking and Biomole Research. Infrastructure.)
 - 250以上の欧州各国のBioBankを統合
- **オランダ Lifeline**
 - 165000人北部オランダ 2006年開始 30年間の追跡、3世代コホート（世界初）
- **米国 Precision Medicine Initiative, Genome Cohort** : “All of Us”コホート
 - これまでのBiobank（例えばBioVUなど）を集めて100万人のゲノムを集める

ビッグデータ医学/医療の2つの流れによる 大規模な生命情報DB/KBの出現と利用

- ヒトゲノム解読計画以降急速に進展
 - Hapmapプロジェクト, 1000 genome, がんICGC, TCGA, TopMED
 - ゲノム変異・多様体
 - dbSNP, HGMD, **Clinvar**, **Clingen**, OMIM, GWAS catalog
 - 表現型との対応: dbGaP, EGA
 - 遺伝子発現プロファイル
 - 疾患特異的transcriptome: **GEO**, **ArrayExpress**,
 - 薬剤特異的transcriptome: **c-Map**, **LINCS**
 - タンパク質
 - 3次元構造: PDB, Swiss-Prot,
 - タンパク質間相互作用: **HPRD**, **STRING**, BIND
 - 分子ネットワーク、パスウェイ
 - KEGG, TRANSFAC, BioCyc, Reactome
- 各種バイオバンク症例ベース（制限アクセス）
 - UK biobank, BMBRI, 東北メディカル・メガバンク
- これらの大規模DB/KBを組合せてゲノム医療/創薬を推進

医学/医療へのビッグデータの衝撃

	HiSeq 2500	HiSeq Proton
本体価格	約1億円	約2500万円
モード / チップ	ハイブリッド / フロント	カブリッド / フロント
解読速度	118	2798
リード長 (bp)	2 x 100	2 x 150
シーケンス長 (G)	8960	8320
設置コスト (ヒト1人ゲノム)	約1万円	約1万円

次世代シーケンサの登場
シーケンス革命 (2007)



コストレスで高精度な網羅的分子情報の出現

1. ゲノム・オミックス医学/医療の進展

— Clinical Sequencingによるゲノム・オミックス医療の臨床実装の急速な進展

2. Biobank/ゲノムコホートの世界的普及

— 個別化医療/予防の情報基盤として普及

3. 大規模な生命情報DB/KBの出現

— ゲノム・オミックスによるDB/KBの膨大化

わが国での現状「ゲノム医療元年 (2016)」

■「ゲノム医学実現推進協議会」(中間報告) 2015.7

研究費を用いた試行的ゲノム医療であるが、いくつかの医療施設でゲノム・オミックス医療が試行されている

●例：がんの網羅的分子診断と個別化治療

- 国立がん研究センター (Top-gear, SCRUM-Japan)
 - ドライバー遺伝子の診断。分子標的薬の治験グループに割当て
 - がんのゲノムパネル：来年先進医療 (7施設)
- 岡大, 京大, 北大, 千葉大 病院併設型BB
- 2018年度より「がん複数遺伝子パネル」先進医療 (7か所) 開始

■AMED (日本医療研究開発法人) がゲノム医療を推進

●IRUD (Initiative on Rare and Undiagnosed Disease)

未診断疾患の原因遺伝子をIRUD拠点病院が審査して解析センターがシーケンシング。その後、DB化する。

- ゲノム医療実現推進プラットフォーム事業
- 臨床ゲノム情報統合DB事業

ゲノム医療では、米国と水を空けられている。しかし、Biobank Genomic Cohortでは我が国の状況はそれほど遅れてはいない。Biobank準拠のゲノム医療/創薬を推進する方針を基本とすべき

医療の「新しいビッグデータの革命性」

～ゲノム・オミックスデータの基軸的な特徴～

＜目的もデータ特性も従来型と違う＞

従来の医療情報の「ビッグデータ」

Big “Small Data” ($n \gg p$)

医療情報・疫学調査では属性数：数十項目程度

— 目的：Population MedicineのBig Data

⇒個別を集めて「集合的法則」を見る

網羅的分子情報などのビッグデータ

Small “Big Data” ($p \gg n$)

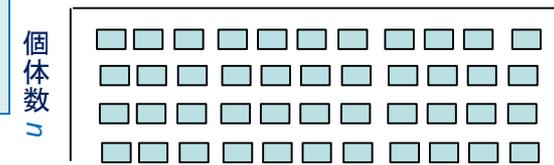
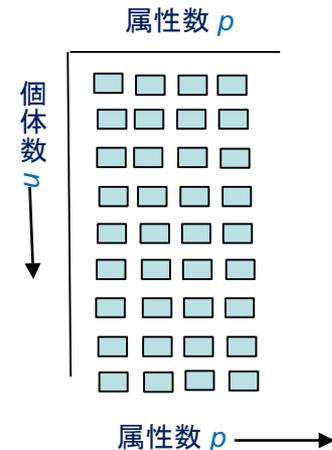
1 個体に関するデータ属性種類数が膨大

属性(p)に比べて個体数(n)少数:従来の統計学が無効

「新 np 問題」：GWASは単変量解析の羅列

— 目的：例えば医療の場合 個別化医療 Personalized Medicine

⇒大量データを集めて「個別化パターン」の多様性を抽出



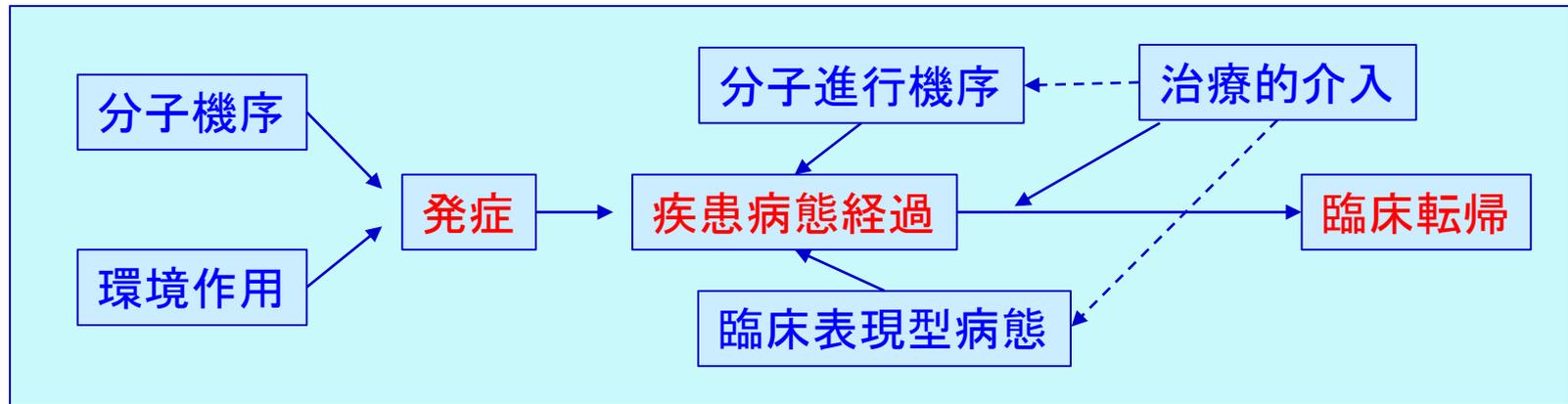
新しいデータ科学の必要性

医療の「ビッグデータ」革命は どんな既存のパラダイムに変革しているか

- Population medicineのパラダイム転換
 - <One size fits for all>のPopulation医療はもはや成り立たない
 - 個別化医療 “Personalized (Precision) medicine”
 - 個別化医療実現のために<個別化・層別化パターン>がどれだけ有るか
網羅的に調べる：どこまでの粒度で個別化・層別化すればよいか
- Clinical research（臨床研究）のパラダイム転換
 - 臨床研究を科学にする従来の範型RCTは、個別化概念を取扱えない
 - <statistical evidence based>呪縛からの解放
 - 「標本」統計・「推測」統計学に制約されない臨床研究
 - Real World Data・ビッグデータからの知識生成（BD2K）

ゲノム・オミックス医療の課題

課題 1 対応する非ゲノム病態データの 検証的「情報化」



病態経過オントロジー

疾患分子発症進行機序（生体分子ネットワーク）

対応する非分子機序の明確化

環境発症要因 臨床表現型情報 治療介入効果

臨床表現型との統合(phenotyping)

臨床表現型データ検証的抽出、非構造化データ障壁

electronic **M**edical **R**ecords + **G**enomics (NHRI-funded)

phase I (2007-2011) EMR-basedゲノム研究の探求

- EMR(臨床phenotyping)とbiorepositoryに基づくGWAS等 (EMR-based GWAS) が可能か。
 - 開始時はGWAS全盛時代。ゲノム医療の臨床実装は始まらず
- 電子カルテより臨床表現型情報抽出 phenotypingルール
- 計画開始時参加施設 : Mayo, Vanderbilt Univ., Marshfield, Univ. Washington, Northwestern Univ.など5施設,

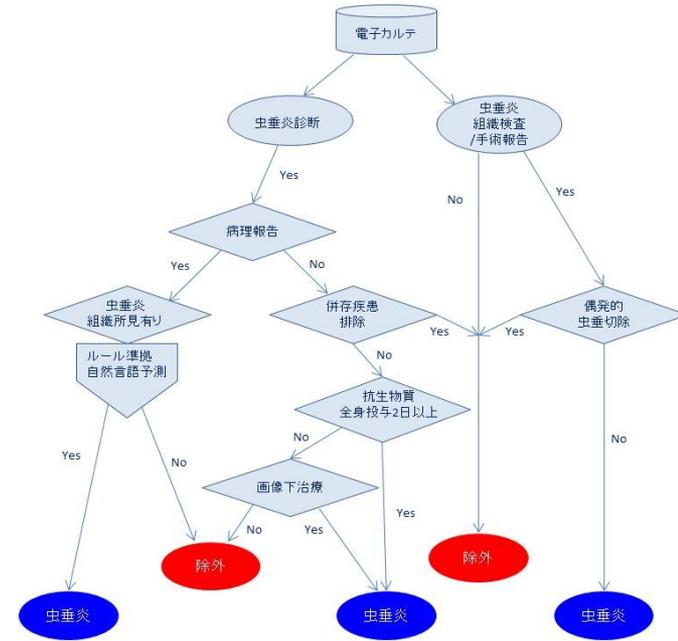
phase II (2011-2015) 臨床実装へ舵を切る

- **MCWの臨床実装のインパクト**、Vanderbiltの先制PG x
- 電子カルテと遺伝情報の統合
 - 電子カルテへのゲノム情報の統合
 - **PheKB** (Phenotype Knowledge Base)
 - ゲノム医療の実装、PGxの臨床応用
 - 結果回付 **Return of Result**, ELSI等
- 4つのサイトが新しく加わる
 - 小児病院グループとMount Sinai, Geisinger

phase III : 2015より始まる

NHGRIのコンソーシアムと連携

- **CSER** “Clinical Sequencing Exploratory Research”



PheKB: phenotyping ルール



課題2 生命医療情報のビッグデータ化による 「革新的(innovative)知識」発見の困難性

- 臨床ゲノム医学
 - 全ゲノム配列の普及、多層オミックス情報の収集、分子画像の発展
 - ビッグデータ化：超多次元相関ネットワーク
 - 〈網羅的分子情報と臨床表現型情報〉の相関
- 予防ゲノム医学
 - バイオバンクの大規模化、国際連携によるバーチャル連携
 - 〈遺伝的素因と環境/生活様式要因〉の相互作用と発症の相関ネットワーク

いずれも超多次元複雑ネットワークの縮約理論

人工知能 Deep Learning への期待

- 機械学習のこれまでの限界

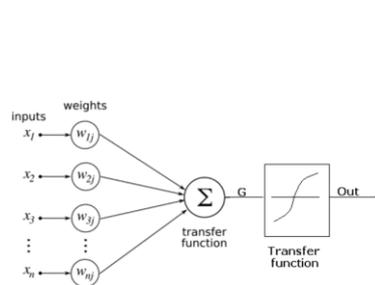
- 「教師あり学習」

- 分類対象の特徴と正解を与え学習機械 (AI) を構築

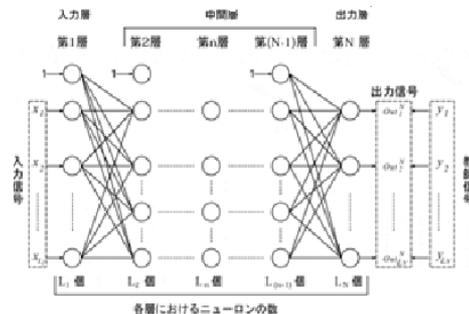
- Deep Learningの革命性

- 「教師なし学習」

- 対象の特徴表現や対象の高次特徴量を自ら学ぶ



神経情報素子

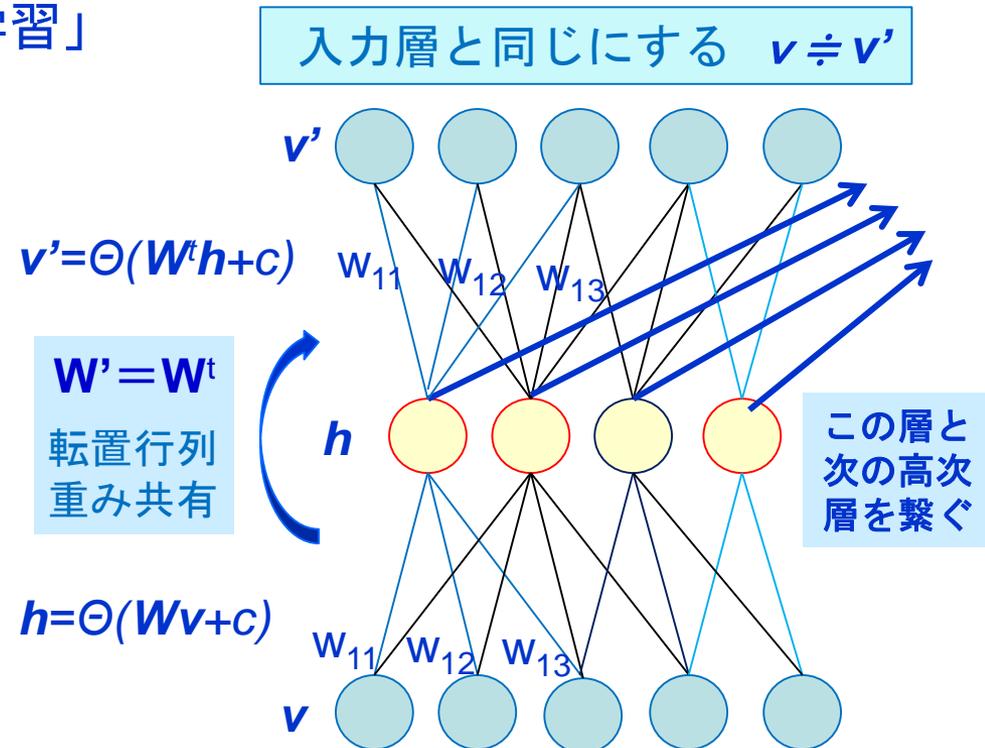
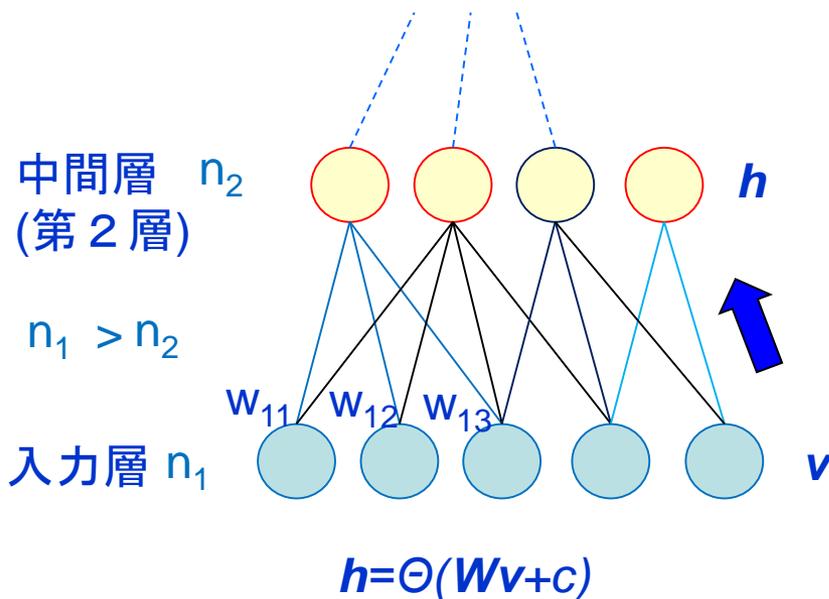


多層ニューロネットワーク



DLの革命点 Autoencoder

- 対象に固有な**内在的特徴**を学ぶ**自己符号化の原理**
- 格段ごとに入力の少ない中間層を入力へ逆投影して復元できるか
- 次元を圧縮され可及的に復元する ($1000_{\text{nodes}} \Rightarrow 100_{\text{nodes}} = ? \Rightarrow 1000_{\text{nodes}}$)
 - できるだけ**復元に効果的な特徴量**を探索する
 - 内在的な特徴量**を見出す
- 最終層で人との対応「教師あり学習」



「ビッグデータ」のData 縮約原理

問題点 属性項目数(p) \gg サンプル数(n)

p : 数億になる場合あり n : 多くても数万、通常数千



これら膨大な属性変数がすべて独立ならばビッグデータの構造解析は不可能。単変量解析の羅列 (GWASのManhattan Plot) しか可能でない



ビッグデータ・スパース仮説

ビッグデータは、多数であるが属性値数より少ない独立成分が基底となって、相互にModificationして構成されている。
(独立成分の推定は、サンプル数とともに増加する)

データ構成性の原理 (principle of compositionality)

Deep Learningによる 多次元ネットワーク縮約法

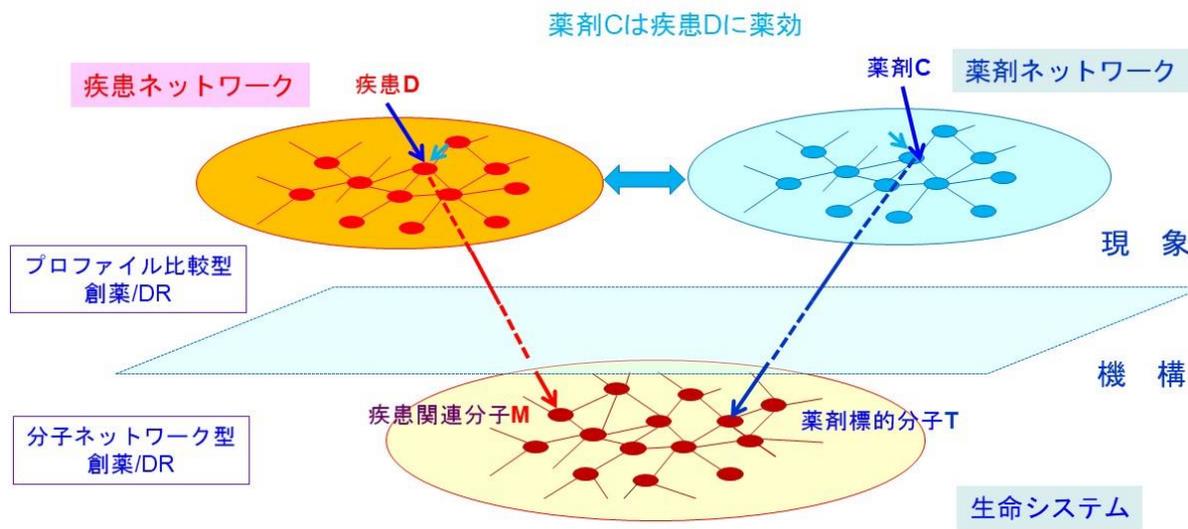
(Hase, Tanaka 2017)

- 医療・創薬ビッグデータへの応用性高い
- 超多次元ネットワーク情報構造の急増
 - ゲノム医療<網羅的分子情報–臨床表現型情報>
 - ゲノムコホート<遺伝素因–環境要因(生活習慣)>
- Deep Learning-based Network Contraction
「DLネットワーク縮約法」
 - 超多次元ネットワーク情報構造⇒
少数の特徴的ネットワーク基底に分解
- 線形分解ではない。非線形分解で基底への射影

人工知能応用としての AI創薬

タンパク質相互作用ネットワークでの 疾患-薬剤-標的分子の関係性の学習

- ビッグデータ創薬/DR
 - タンパク質相互作用ネットワーク上での有効性予測
 - 基準指標：「疾患関連分子」と「薬剤標的分子」の距離
 - 判定情報量が不足
- AI創薬/DR
 - ビッグデータ創薬/DRの限界をAI学習で補完
 - 既成の疾患-薬剤-標的分子の正例を学習（DrugBank）
 - 疾患関連分子と標的分子のPPIN位置関係性をDLで学習
 - 学習された関係性より各分子の**標的分子の有効性を判定**

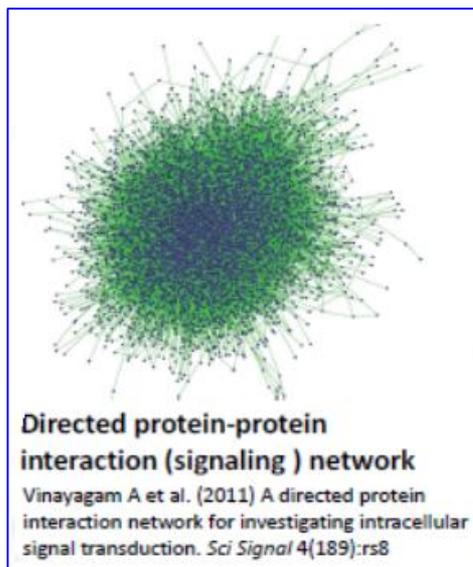


特徴的ネットワーク基底への分解

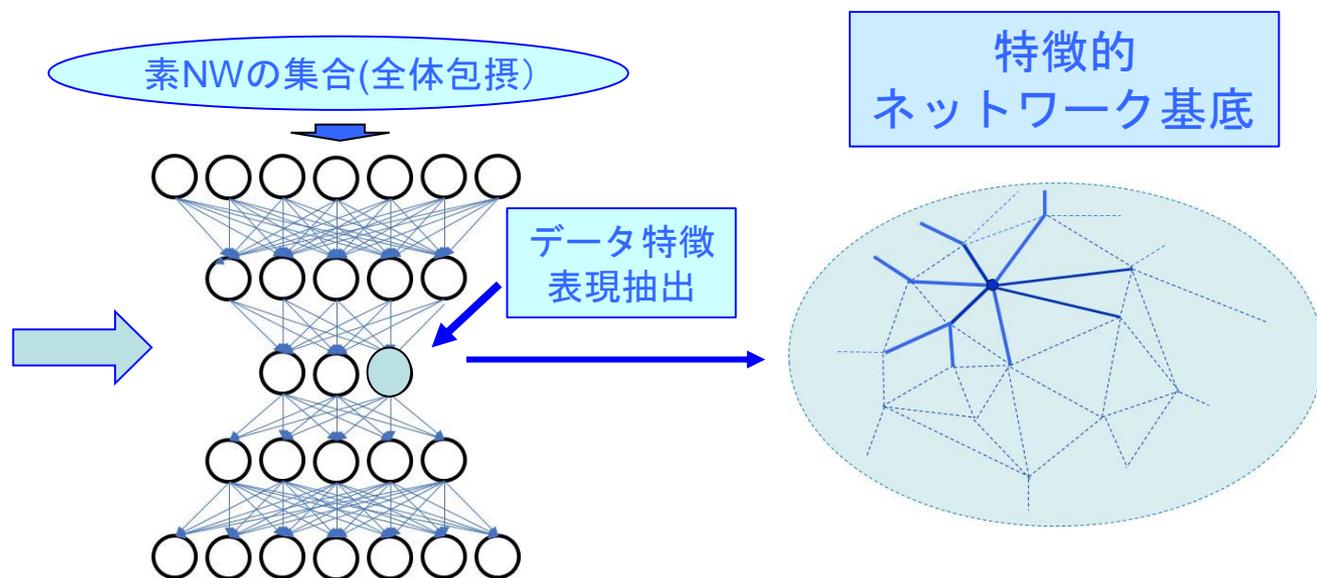
特徴的ネットワーク基底の和に縮約

特定のノードを起点とした素NW（部分NW）の集合
全体NWを包摂する集合にDL反復自己学習

特徴的ネットワーク基底：トポロジーのみの構造/頻度構造



PPIネットワーク



Deep Learningによる創薬・DR

1) 生体ネットワーク (PPIN) 特徴量の抽出

- タンパク質相互作用ネットワーク(PPIN)のNW結合を学習し**特徴表現** (特徴NW基底) を出力。
- 学習集合を部分ネットワークの集合から決める
- ノードを起点とした素NWでPPIN全体を覆う集合

2) 多層Deep Auto-encoderのDLで学習.

- 特徴的NW基底の「教師無し」学習
- 次元縮約による特徴的NW基底の抽出

3) DL特徴NW基底空間における正例補完

- DrugBankからの正例とその増加 (SMOTE法)

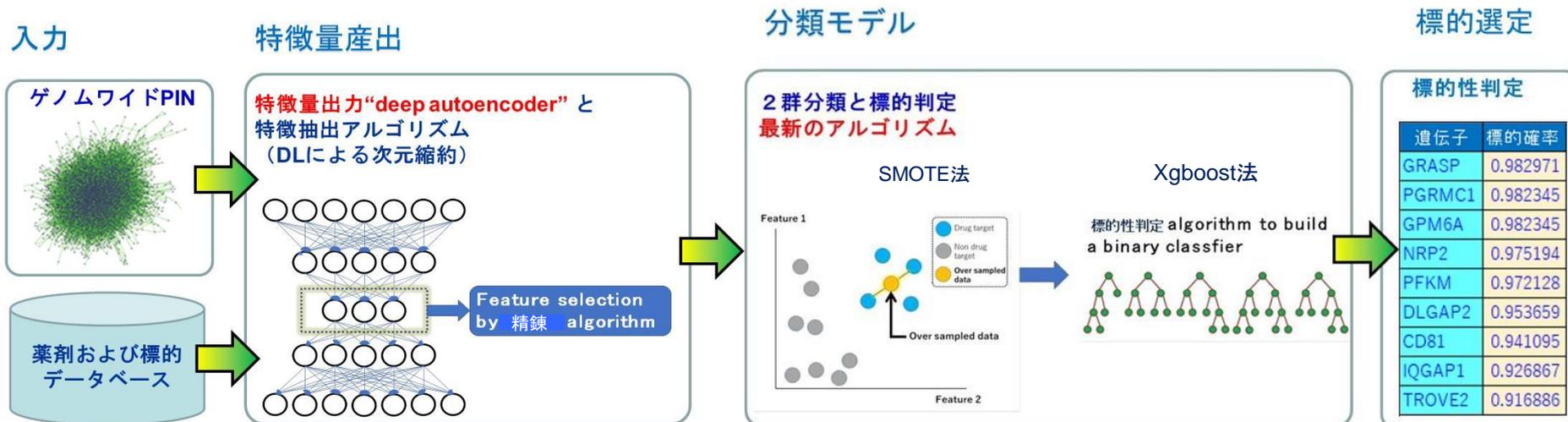
4) DL特徴NW基底量を用いた機械学習分類

- Xgboot法などを用いたDL特徴量からの判別ネットワーク・タンパク質の標的性の判定

Deep Learningによる創薬・DR

分類部 DrugBankを利用した 当該分子を標的とする既製薬剤の探索

既製薬剤がない→新規薬剤探求（創薬）
既製薬剤がある→DRの検討



従来の機械学習（Random Forrest）と同じ成果は得られている

実験的研究との付合

PGCM1 : progesterone receptor membrane 1

Journal of Neurochemistry
JNC

JOURNAL OF NEUROCHEMISTRY | 2017 | 140 | 561-575 | doi: 10.1111/jnc.13917

ORIGINAL ARTICLE

Small molecule modulator of sigma 2 receptor is neuroprotective and reduces cognitive deficits and neuroinflammation in experimental models of Alzheimer's disease

GPM6A : Glycoprotein M6A

INTERNATIONAL JOURNAL OF MOLECULAR MEDICINE 25: 447-455, 2010

Characterization of changes in global gene expression in the brain of neuron-specific enolase/human Tau23 transgenic mice in response to overexpression of Tau protein

CD81:Tetraspanins family

frontiers in Molecular Neuroscience

MINI REVIEW published: 21 December 2016 doi: 10.3389/fnmol.2016.00149

The Emerging Role of Tetraspanins in the Proteolytic Processing of the Amyloid Precursor Protein

Lisa Seipold and Paul Saftig*

Institut für Biochemie, Christian-Albrechts-Universität zu Kiel (CAU), Kiel, Germany

OPEN ACCESS Freely available online

PLOS ONE

Alzheimer's Therapeutics Targeting Amyloid Beta 1-42 Oligomers II: Sigma-2/PGRMC1 Receptors Mediate Abeta 42 Oligomer Binding and Synaptotoxicity

Nicholas J. Izzo¹, Jinbin Xu², Chenbo Zeng², Molly J. Kirk^{5,9}, Kelsie Mozzoni¹, Colleen Silky¹, Courtney Rehak¹, Raymond Yurko¹, Gary Look¹, Gilbert Rishton¹, Hank Safferstein¹, Carlos Cruchaga⁶, Alison Goate⁶, Michael A. Cahill¹⁰, Ottavio Arancio⁷, Robert H. Mach², Rolf Craven⁴, Elizabeth Head⁴, Harry Levine III³, Tara L. Spires-Jones^{5,8}, Susan M. Catalano^{1*}

DLGAP2 : DLG-Associated Protein 2

Journal of Alzheimer's Disease 44 (2017) 181-194
DOI: 10.1007/s12064-016-0420-8
Kim Park

Genetic Variation in Imprinted Genes is Associated with Risk of Late-Onset Alzheimer's Disease

PFKM: Phosphofruktokinase

Cytotechnology (2016) 68:2567-2578
DOI 10.1007/s10616-016-9980-3

ORIGINAL ARTICLE

Neuroprotective effect of Picholine virgin olive oil and its hydroxycinnamic acids component against β -amyloid-induced toxicity in SH-SY5Y neurotypic cells



第2世代の ゲノム・オミックス医療

ゲノム医療の第2世代

成功した臨床実装

1. 希少先天遺伝疾患の原因遺伝子を病院の現場でシーケンサにより同定
2. がんのドライバー遺伝子変異を同定、適切な分子標的薬を処方
3. 患者の薬剤の代謝酵素の多型性を先制的に同定し、副作用を防ぐ

しかし

多因子疾患の機序/発症予測は無着手である

- 「単一遺伝的原因」 帰着アプローチの限界
- 「行方不明の遺伝力」の主要な原因
複数の疾患関連遺伝子間の相互作用: $G \times G$
環境と遺伝子の相互作用が: $G \times E$

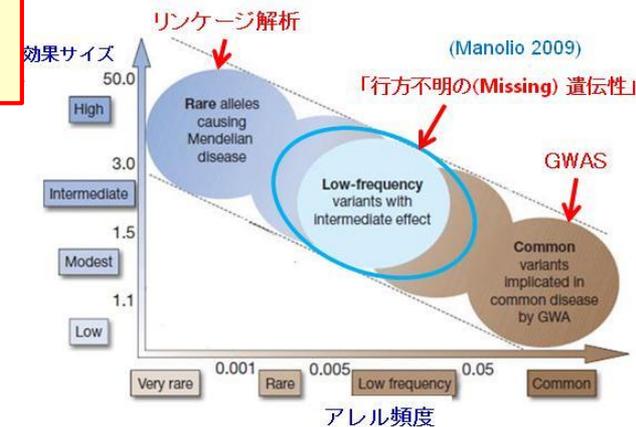
SNPの相対リスク
低い(1.1~1.3)理由
 $G \times E$ 組合せ特異的効果
を環境要因の平均



多因子疾患は個人の<遺伝的体質と環境要因>の
<相互作用の結果。シーケンスだけでは解明不能

疾患発症の遺伝要因と環境要因の相互作用は
加算的 ($G \oplus E$) でもなく乗算的 ($G \otimes E$) でもない
< (G, E) 組合せ特異的な効果 > である

例 大腸がんの遺伝要因と環境 (生活習慣) 要因



大半の疾患の基礎としての 「遺伝素因X環境要因」の相互作用

一部の単一遺伝病を除き、大半の疾患
(Common diseases)の発症は

疾患発症の相対リスク=

遺伝要因(G:genome) X 環境要因(E:exposome)

相互作用は加算的でもなく乗算的でもない

<(G,E) 組合せ特異的な効果>である

GWASでSNPの相対リスクが低い
(1.1~1.3)理由: GxE組合せ特異
的效果を環境要因の全てに亘って
平均しているからである



発達プログラム説 DOHaD

(Developmental Origin of Health and Disease)

- オランダ飢饉
 - 第2次大戦末期、ナチスの封鎖、約半年間酷い飢饉
 - 飢饉の期間に胎児、戦後30年
 - 成人期:肥満,糖尿病,心筋梗塞,統合失調
- Baker仮説：英国心筋梗塞増加
- エピジェネティック機構
 - 過度な低栄養：肝臓のPPAR α/γ （儉約遺伝子）メチル化低下・遺伝子発現がオン
 - エピジェネティック変化は可変：短期的変化、長期的「記憶」次の世代も



オランダ
飢饉 (1944)

環境因子

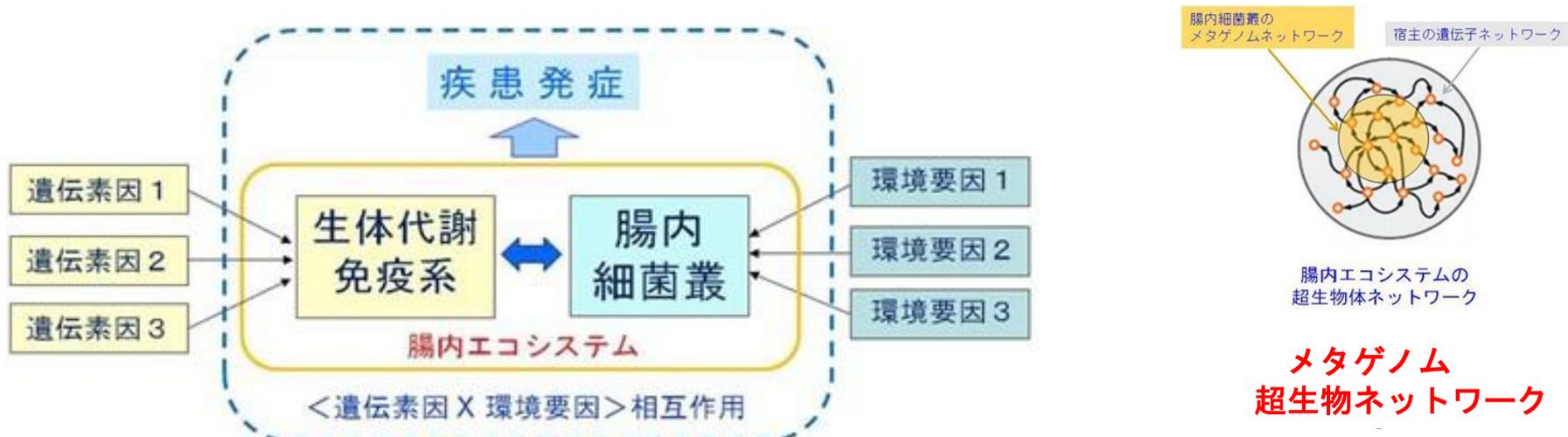
Epigenome変化

遺伝子発現調節

疾病発症

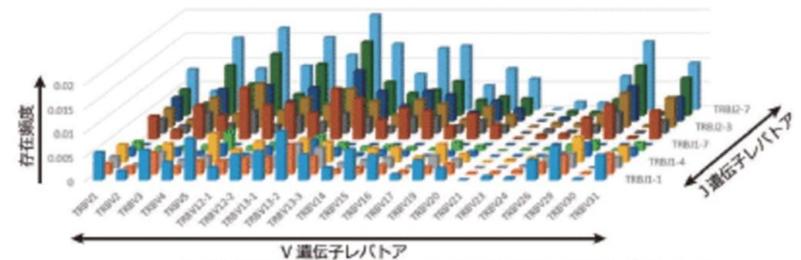
腸内細菌叢microbiome：メタゲノム

- 疾患の環境発症要因（exposome）
 - 腸内microbiome：環境要因の最大の1つ
- 腸管微生物叢（gut microbiome）
 - 約1000種類、100兆個、総重量1～1.5kg, 「**実質的な臓器**」
 - 遺伝子数個人あたり約**50万遺伝子**、総数：数100万遺伝子
- **免疫系、炎症系、粘膜免疫細胞群との相互作用**
 - 食物の難消化性の食物繊維：腸内細菌によって嫌氣的に代謝、酪酸などの「**短鎖脂肪酸**」がエネルギー源となる
 - 食事・栄養物質による環境要因は、腸内細菌叢の代謝物（短鎖脂肪酸やTMAOなど）から宿主の生体機構に相互作用



免疫ゲノム

- 可変領域や相補性決定領域（特にCDR3）のDNAやRNAを次世代シーケンサ(HTS)で解析
- レパトア解析
 - 抗原受容体全体のプロファイルを俯瞰的に把握できる
 - V(D)Jなどの成分を基軸として3次元表示可能。
 - 疾患罹患とともに瞬時に全体像が変化する。
 - 網羅的病態全体像を提示する
 - VDJの使用頻度
 - 多様性(diversity)の変化
 - 疾病/加齢レパトア分布変化
- 臨床シーケンスに含まれる
- 3次元分布の特徴分析



(レパトア・ジェネシス社)

第2世代のゲノム・オミックス医療

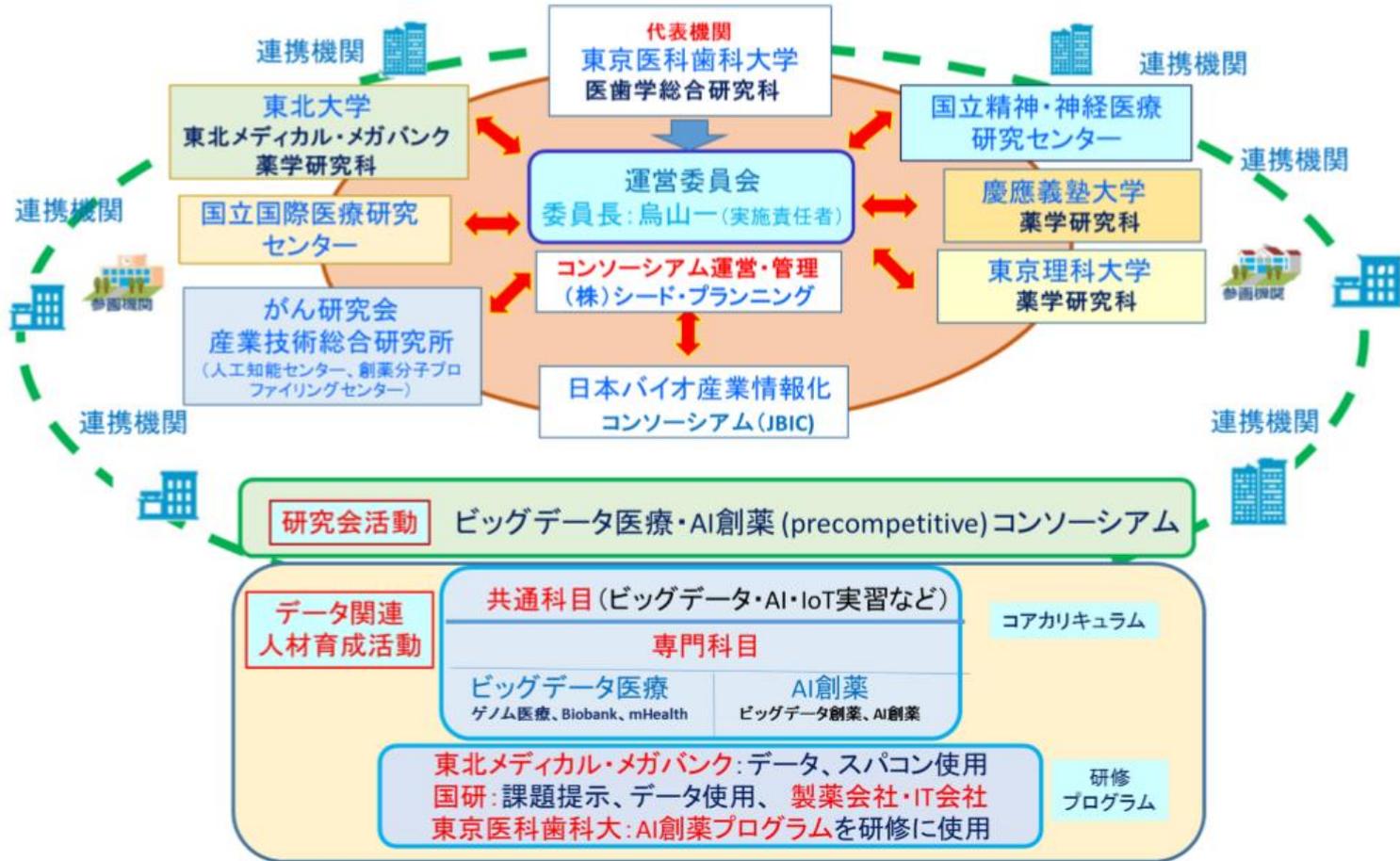
- 生涯的全体性においてその個人の疾患可能性の全体性を把握し、個別化予防、個別化治療に取り組む
- ゲノム・オミックス情報と医療・健康
 - **Clinical Sequencing**のインパクト
- **第1世代ゲノム医療**
 - ゲノムの変異・多型性の個別性に基づく
- **第2世代のゲノム医療**
 - 多因子疾患が対象、環境情報との相互作用
 - エピゲノム、メタゲノム・免疫ゲノムなど

疾患メタ・オミックス修飾

学会の活動

- 日本オミックス医療学会シンポジウム『「第二世代のゲノム医療実現に向けて」～医療におけるビッグデータ最前線～』
2017年6月21日(水) 東京
- 日本オミックス医療学会シンポジウム「AI医療・AI創薬」2016年12月14日(水) 東京
- 日本オミックス医療学会 シンポジウム「免疫ゲノム～免疫ゲノムの多様性から疾患への介入を探る」2016年10月18日 東京
- 日本オミックス医療学会 シンポジウム「microbiome創薬」
2016年6月24日 東京
- 日本オミックス医療学会シンポジウム2016年3月11日「ゲノム医療元年」東京
- 日本オミックス医療学会シンポジウム「AI創薬」2015年12月14日
- 日本オミックス医療学会 シンポジウム「ビッグデータを創薬・育薬へ」2015年9月16日

ビッグデータ医療・AI創薬コンソーシアム



2018.2.23-25, Harvard/MIT/TMDU - Datathon

ご清聴ありがとうございました

