

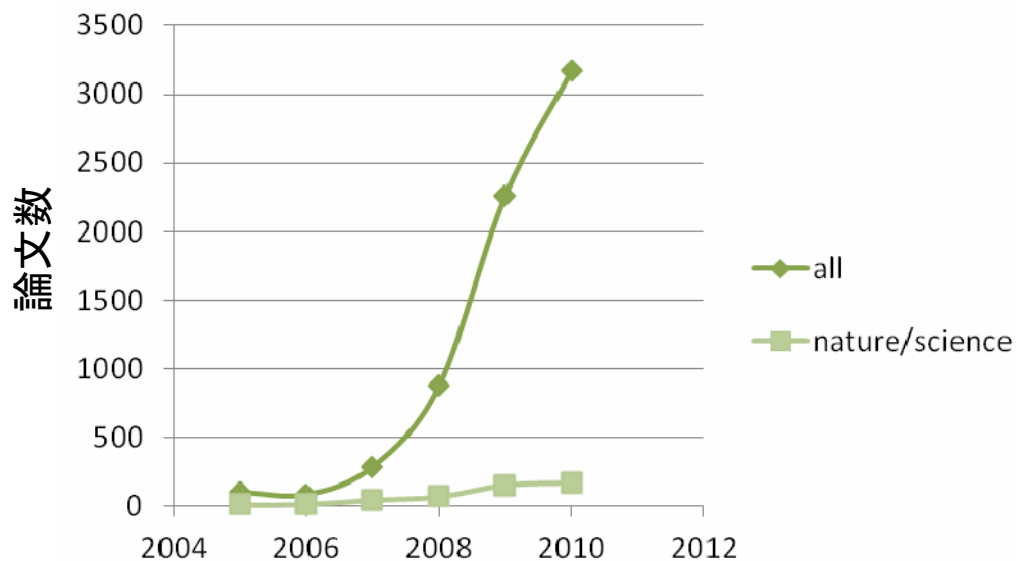
統合の日シンポジウム
October 5, 2011

ヒトゲノムバリエーション データベースの開発

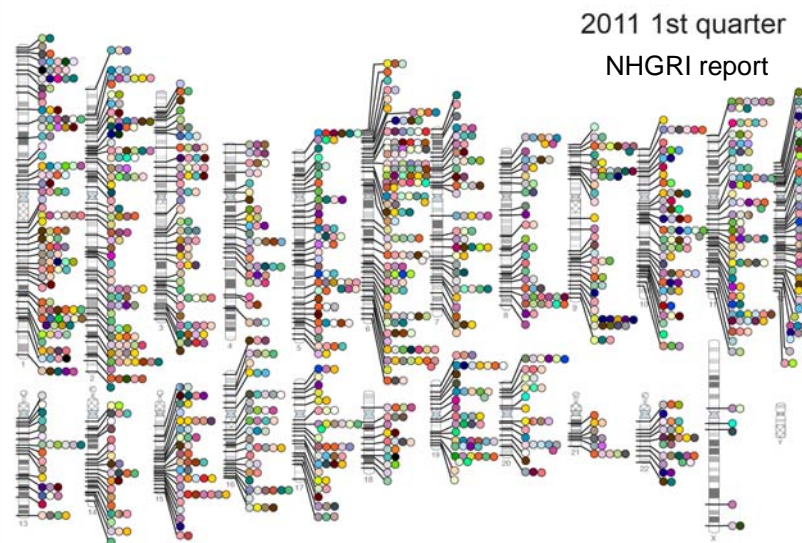
東京大学大学院医学系研究科
人類遺伝学分野
徳永勝士

背景 1 – ゲノムワイド関連解析 (GWAS)

ヒトゲノムの解読完了、HapMap PJの完了とSNPタイピング技術の飛躍的向上により大規模なゲノムワイドな関連解析が多数進行



論文数の推移

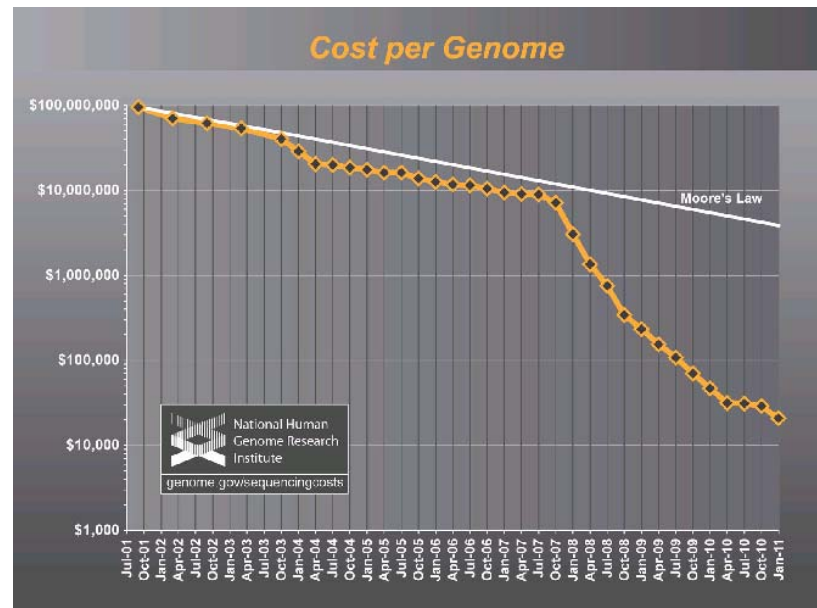


Published Genome-Wide Associations (03/2011)
1,319 published GWA at $p \leq 5 \times 10^{-8}$ for 221 traits

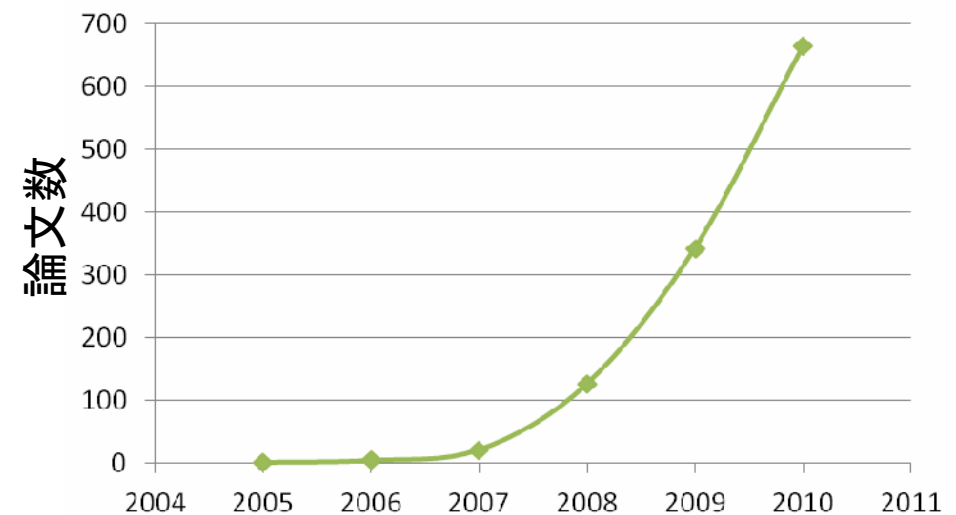
1000人、100万SNP解析の生データ: 10TB程度

背景 2 – 次世代シーケンサー

次世代シーケンサーの技術革新、ゲノム解読コストの低下により、whole genome/exome sequencingによる疾患変異の探索も多数進行



ゲノム解析コストの推移

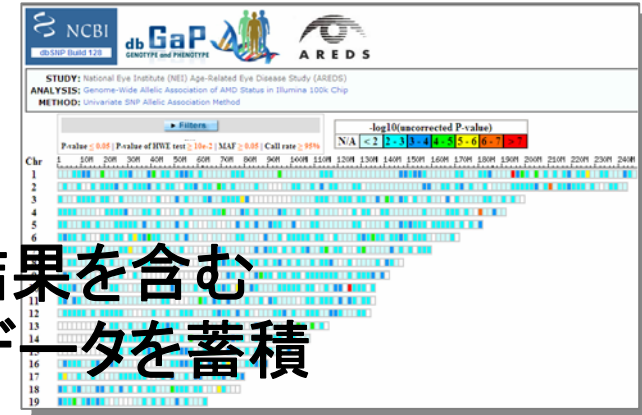


論文数の推移

背景 3 — 海外の関連するデータベース

1) NCBI

- dbSNP SNP情報を蓄積
- dbVAR 構造多型のデータを蓄積
- dbGaP GWAS, 次世代シーケンサー結果を含む genotype-phenotype に関するデータを蓄積



2) EBI

- EGA GWAS, 次世代シーケンサー結果を含む genotype-phenotype に関するデータを蓄積



3) HGVS (human genome variation society)

- LSDB (locus specific database) のリンク集

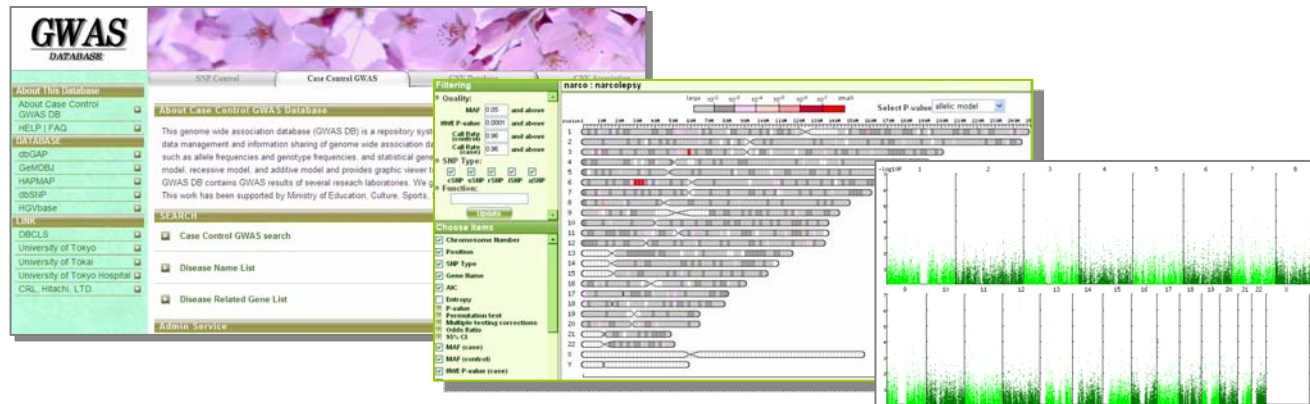
4) その他、プロジェクト単位のDB

(1000 genome PJ, 国際がんゲノムコンソーシアム等)

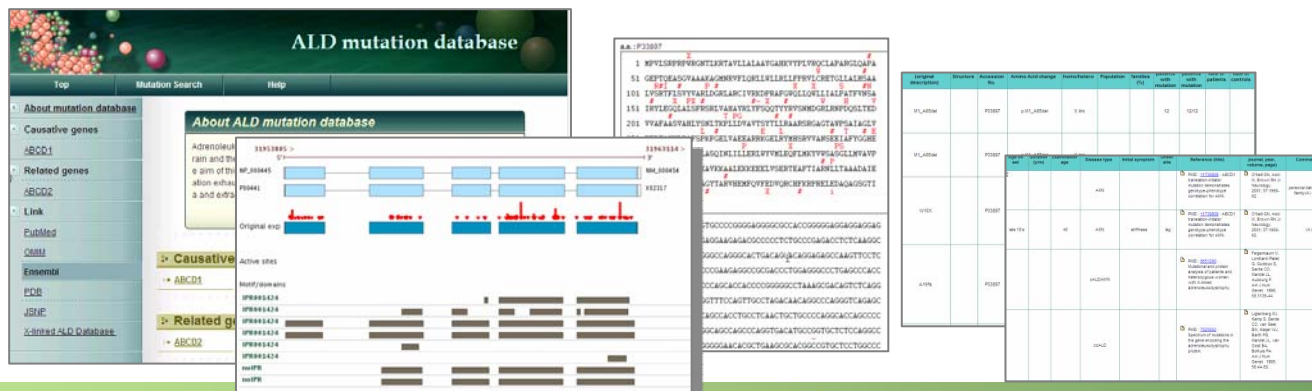
統合データベースプロジェクト(～H23.3)での取り組み

東京大・医:徳永勝士、辻省次、 国立遺伝研:井ノ上逸朗、 日立中央研:小池麻子
(<https://gwas.lifesciencedb.jp/index.html>)

GWAS-DB GWAS結果のDB化、研究者間の情報共有を促進



Mutation DB 神経変性疾患の臨床情報と変異情報の関係を俯瞰



ヒトゲノムバリエーションデータベース (H23.4~)

高速大量のSNPタイピング技術と次世代ゲノム配列解析技術の発展
疾患関連遺伝子多型・変異の探索が、世界レベルで大規模に進行
ヒトゲノムの多様性および疾患関連遺伝子・変異情報が急激に増加

1 検体あたりゲノムの0.1%にあたる約300万のSNV(single nucleotide variant)が検出され、その10%程度はdbSNPに未登録

疾患発症の機序は複雑

- 1) 複数因子が複雑に疾患に関与
- 2) 同一遺伝子変異の複数疾患への関与
- 3) 同一疾患における変異部位・種類特異的症状の存在
- 4) 原因・関連遺伝子(変異)の集団間差異の存在

多型・変異の医学生物学的意義の理解には、個々の集団における頻度情報が必要不可欠



疾患・変異・臨床情報の関係を整理・体系化し、得られた成果・情報を公開・共有することにより、疾患機序の解明や個別化医療の実現に貢献

目的と構想

目的:

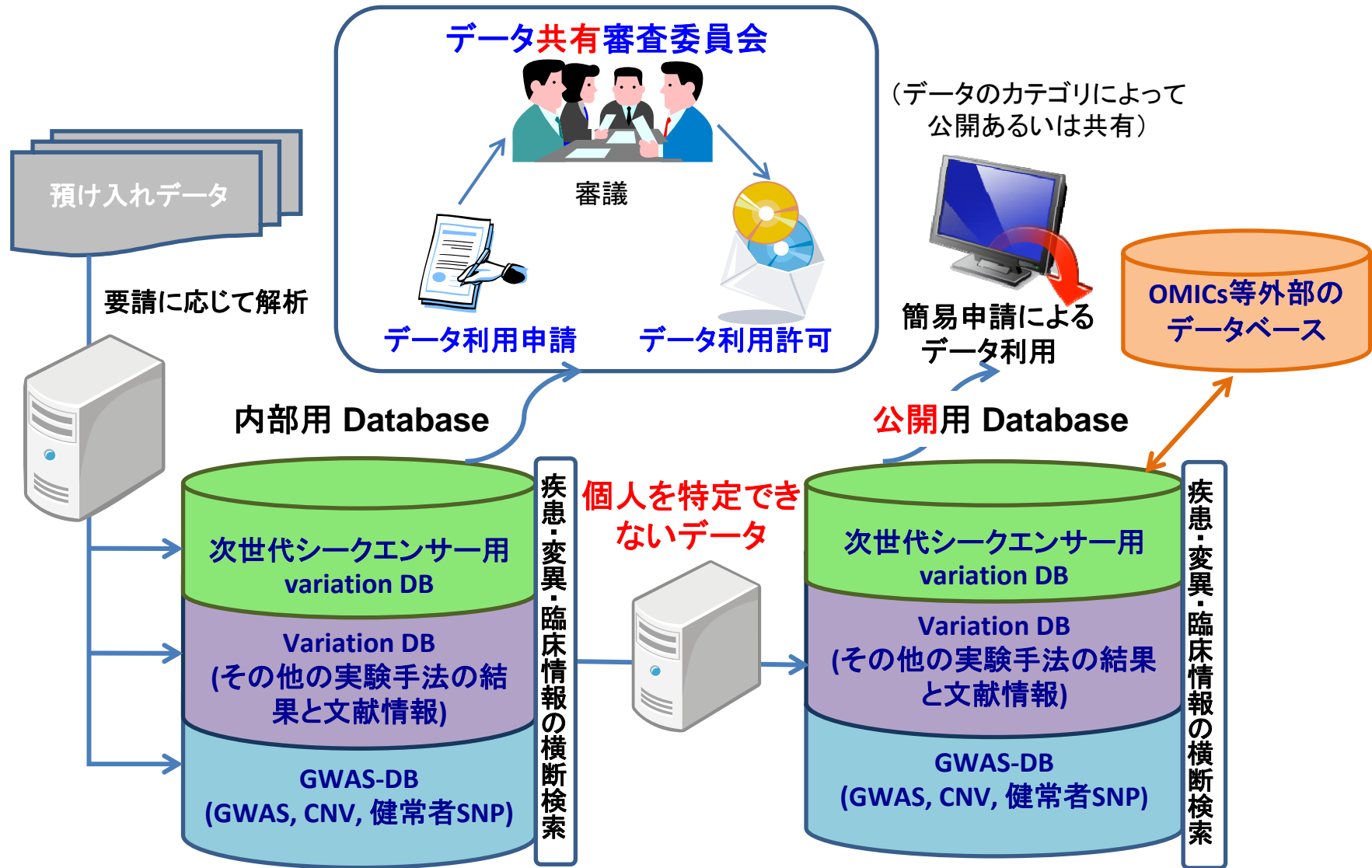
変異・疾患・臨床情報を整理・体系化し、成果・情報を俯瞰可能とすると共に、健常者のゲノム多様性情報を提供する

構想:

- 1) 次世代シーケンサー、その他の解析法(SNP-GWAS、CNVなど)によって発見される多型・変異-疾患情報の受け入れ・保管
- 2) 文献情報など過去に報告された疾患感受性、薬剤応答性などに関わる多型・変異データの収集
- 3) 上記データを整理体系化したDBの構築、データの公開と共有
(疾患→多型・変異、多型・変異→疾患を横断的に探索可能)
- 4) 健常者データ: 新規データ、1000人ゲノム、GWAS健常者データなどを用い、SNP/SNV、repeat変異、in/del、CNVなど各種多型・変異のアリル頻度、ハプロタイプ頻度を計算・公開

→ 効率的な疾患遺伝子の探索に役立てる

データベースシステムの概要



GWAS-DBのコンテンツ

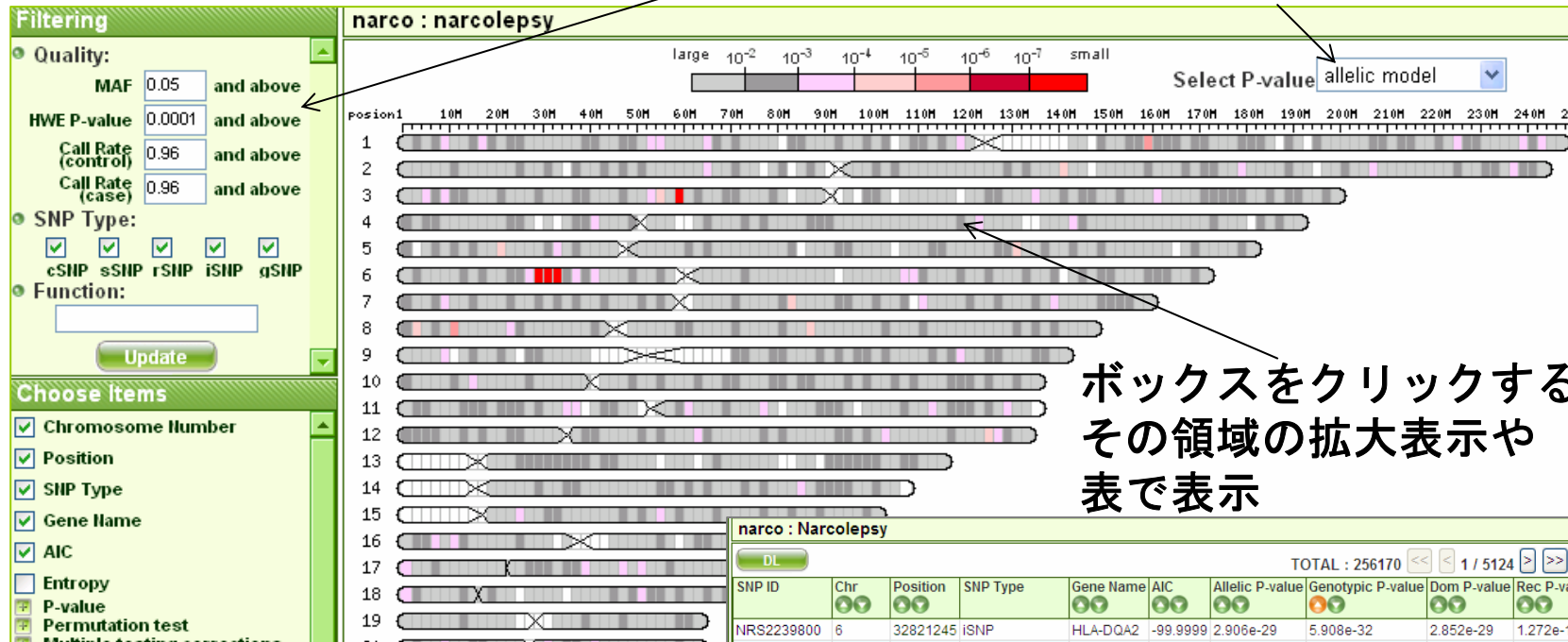
- 30-100万SNPの遺伝子型頻度、アレル頻度、ハーディー・ワインベルク平衡検定値、Call rate等
- genotypic test, allelic test, additive risk model, recessive model, dominant model のP-value, OR, 95% CI, AICなどの遺伝統計値 (plink)
- ハプロタイプもしくはSNPの組み合わせに関する疾患関連性の統計値 (Haploview)
- SNPのアノテーション (機能、染色体上位置、同義/非同義置換など)

GWAS-DB 俯瞰図

<https://gwas.lifesciencedb.jp/index.html>

閾値を変えて表示可能

モデルを変えて表示可能

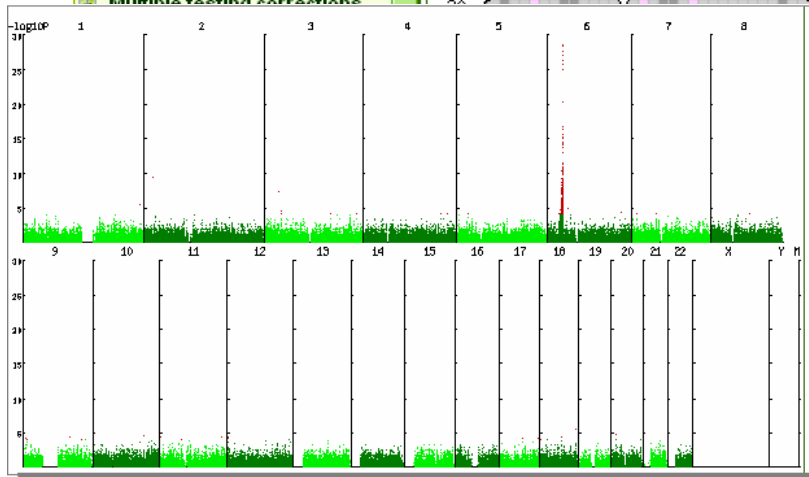


ボックスをクリックするとその領域の拡大表示や表で表示

narco : Narcolepsy

DL TOTAL : 256170 1 / 5124 Reset

SNP ID	Chr	Position	SNP Type	Gene Name	AIC	Allelic P-value	Genotypic P-value	Dom P-value	Rec P-value	Add P-value
NRS2239800	6	32821245	ISNP	HLA-DQA2	-99.9999	2.906e-29	5.908e-32	2.852e-29	1.272e-13	4.218e-2
NRS544358	6	32381136	ISNP	C6orf10	-99.9999	4.158e-29	1.956e-30	1.052e-26	1.007e-14	2.737e-2
NRS926591	6	32413668	ISNP	C6orf10	-99.9999	2.378e-28	1.08e-29	1.145e-25	1.109e-14	9.111e-2
NRS8192590	6	32295761	ISNP	NOTCH4	-99.9999	5.3e-29	1.757e-29	8.585e-28	9.214e-09	2.617e-2
NRS539703	6	32396440	ISNP	C6orf10	-99.9999	7.105e-28	3.572e-29	1.167e-25	3.674e-14	3.164e-2
NRS4959093	6	32421075	ISNP	C6orf10	-99.9999	4.182e-27	1.175e-28	7.841e-25	3.539e-14	6.799e-2
NRS574710	6	32396168	ISNP	C6orf10	-99.9999	3.949e-27	3.397e-28	3.467e-24	2.44e-14	1.449e-2
NRS910050	6	32423632	ISNP	C6orf10	-99.9999	2.557e-26	2.817e-27	2.983e-16	4.39e-22	1.329e-2
NRS9268302	6	32432795	ISNP	C6orf10	-99.9999	1.324e-25	1.511e-26	3.748e-15	9.344e-22	9.149e-2
NRS9268402	6	32449331	ISNP	C6orf10	-92.5310	2.082e-17	4.598e-21	1.54e-12	5.002e-16	5.752e-1
NRS9275134	6	32758590	ISNP	C6orf10	-91.2270	4.256e-17	8.215e-21	7.098e-13	4.525e-15	4.765e-1
NRS241400	6	32979514	ISNP	C6orf10	-86.7120	6.25e-21	1.634e-20	2.188e-16	5.645e-13	2.3e-20
NRS2856688	6	32762618	ISNP	C6orf10	-87.1350	2.767e-16	6.547e-20	1.54e-12	3.595e-14	2.303e-2
NRS3806156	6	32481676	ISNP	BTNL2	-80.5290	1.176e-15	6.219e-19	3.516e-10	9.328e-16	2.19e-16
NRS9268645	6	32516505	ISNP	HLA-DRA	-81.2470	3.466e-15	1.264e-18	4.608e-09	2.711e-16	1.472e-1
NRS9268560	6	32497490	ISNP	C6orf10	-77.7950	8.952e-14	6.878e-18	2.194e-08	7.881e-16	5.953e-1
NRS2076533	6	32471505	ISNP	BTNL2	-77.5070	8.94e-15	8.12e-18	3.465e-10	3.909e-14	5.788e-1
NRS9461799	6	32797507	ISNP	C6orf10	-69.4640	1.065e-15	1.854e-16	1.552e-10	5.968e-13	1.314e-1



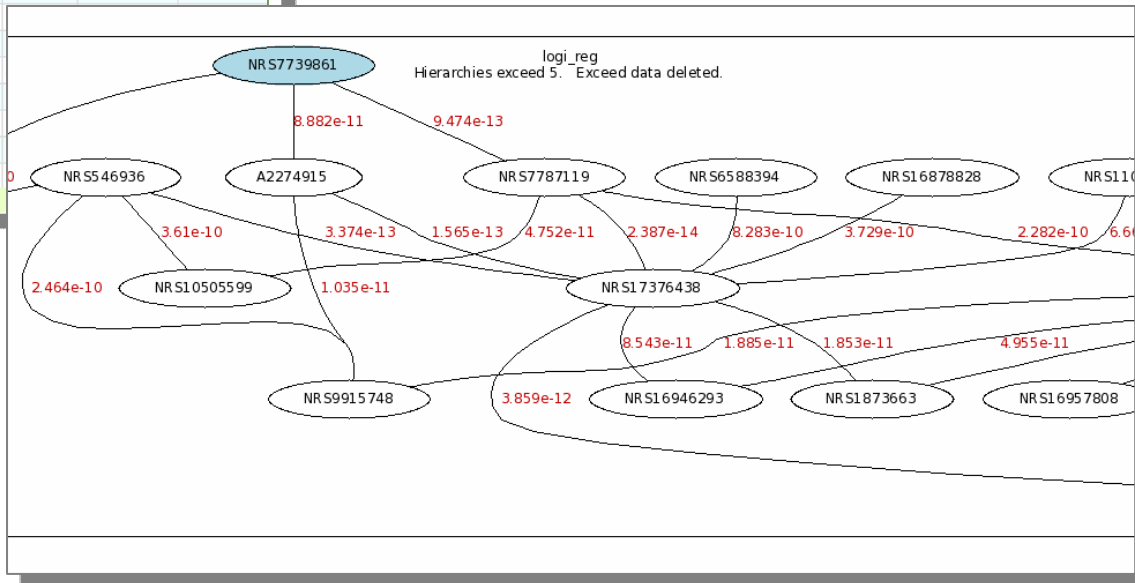
マンハッタン図

Koike A et al. J Hum Genet (2009)

Epistasisの表示

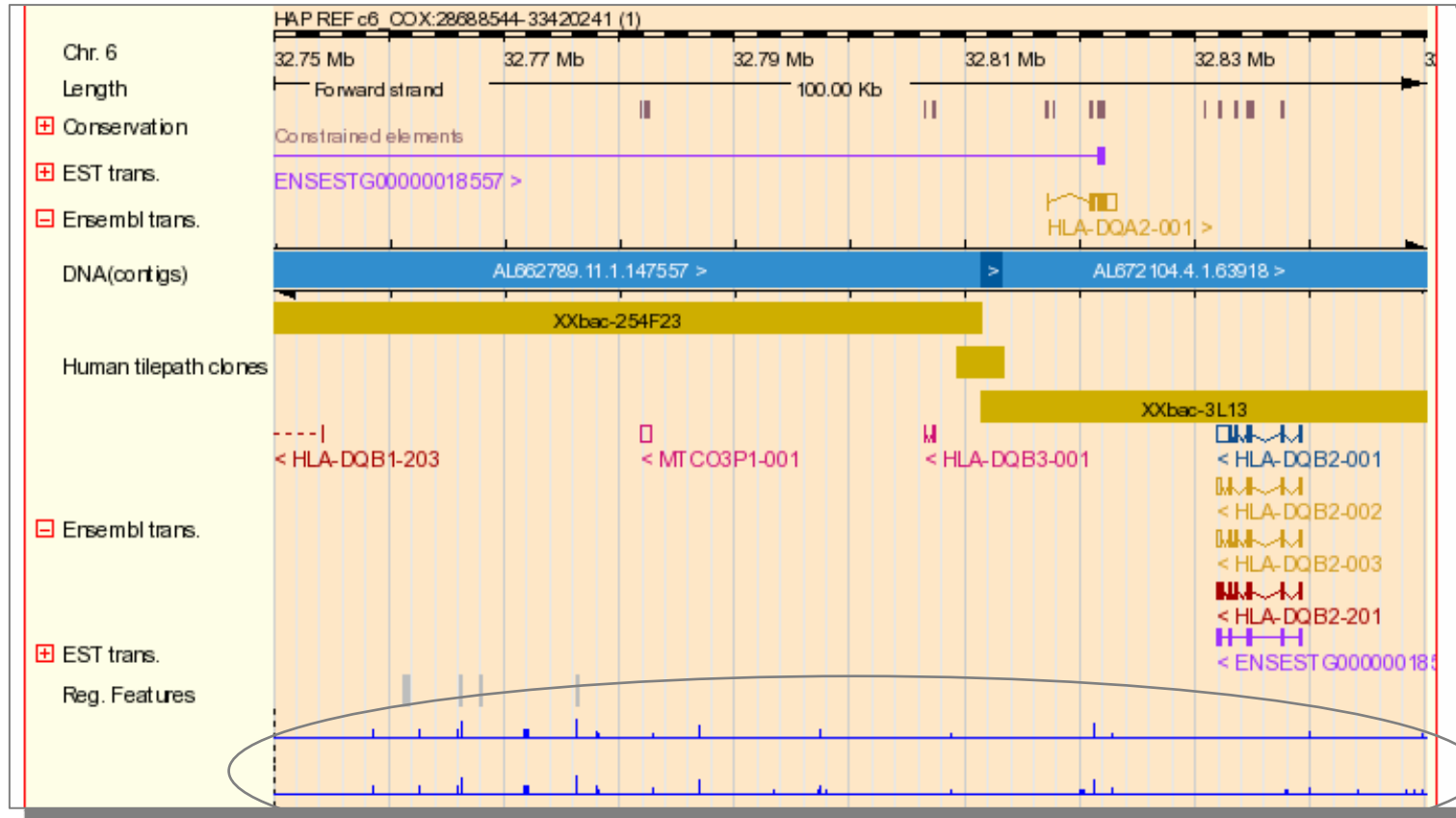
Int Kind	Gene A/SNP-A	Chr	Pos	Gene B/SNP-B	Chr	Pos	Weight
logi_reg	NRS7787119	7	43687934	NRS17376438	10	77725266	2.387e-14
logi_reg	A2274915	9	14250138	NRS17376438	10	77725266	1.565e-13
logi_reg	NRS546936	6	44757114	NRS17376438	10	77725266	3.374e-13
logi_reg	NRS7739861	6	667716	NRS7787119	7	43687934	9.474e-13
logi_reg	NRS17824132	8	73317985	NRS16957808	13	100303160	1.546e-12
logi_reg	NRS17824132	8	73317985	NRS3902307	12	100045651	3.694e-12
logi_reg	NRS17376438	10	77725266	NRS6035066	20	18464588	3.859e-12
logi_reg	NRS7787119	7	43687934	NRS6008536	22	46842252	4.935e-12
logi_reg	NRS7787119	7	43687934	NRS6085297	20	5791479	5.276e-12
logi_reg	NRS1945652	11	80905066	NRS16957808	13	100303160	9.108e-12
logi_reg	NRS17824132	8	73317985	NRS17102227	14	64391947	9.725e-12
logi_reg	A2274915	9	14250138	NRS9915748	17	51518243	1.035e-11
logi_reg	NRS7787119	7	43687934	NRS9915748	17	51518243	1.167e-11
logi_reg	NRS1945652	11	80905066	NRS11090824	22	47038271	1.305e-11
logi_reg	NRS1945652	11	80905066	NRS17102227			
logi_reg	NRS17824132	8	73317985	NRS16956024			
logi_reg	NRS4682115	3	113923788	NRS17824132			
logi_reg	NRS17376438	10	77725266	NRS1873663			
logi_reg	NRS17824132	8	73317985	NRS9915748			
logi_reg	NRS7787119	7	43687934	NRS16956024			
logi_reg	NRS4700811	5	178966007	NRS17824132			

Logistic regressionをはじめとした
Epistasisの登録、表示



DAS server としての機能

他のデータベースからの呼び出しが可能



Entry point:

<https://gwas.lifesciencedb.jp/cgi-bin/das/dsn>

CNV control DB, CNV case-control DB コンテンツ

- 収集対象: GWASで得られたデータを用いたCNV解析
- DBの種類:
 - CNV control DB: 健常者のCNV DB
 - CNV case control DB: CNVのcase-control 関連解析
- 表示
 - 開始、終止位置が微妙に異なるCNVも多く検出されるため、類似CNVをクラスタリングしてマージし、CNVのパターン一覧を見やすくしたオプションも用意
 - 現状の実験精度を考慮し、複数の計算結果を表示

Human Variation DBのコンテンツ

- 収集対象の変異の種類
mutation, small/long insertion/deletion, structural variation, CNV, repeat variation
- 収集対象
 - ・NGSによって検出した健常者/疾患関連変異、実験条件、計算条件
 - ・その他の実験によって検出した変異、実験条件
 - ・文献中に記載される変異、実験条件
- 登録していただくデータ、1000人ゲノムPJのデータ、その他の公開日本人健常者データを元に、まずはconsensus sequenceを作成して、変異情報とともに公開する予定

Human Variation DBのプロトタイプ

Human Variation DB

Filter by Conditions

- Disease
 - RA
 - IDDM
- Experiments
 - Taqman
 - DNA sequencing
 - DNA array
 - Allele specific PCR
 - Taqman 5'-allele discrimination

Update

Chromosome 1 Region 17662588 - 17662737 Show 150bp

GGGCCAGGCTGGGTGCCCAACCCCGACCCACCAACCTCTCCTTACTTIGATGGGATTTGAG AA ATC TCG TTG TGC GCA GAC ATC ACC CGC ACC GGC AAA GTG AAG CCA AC
 CCCGGTCCGACCCACGGGGTTGGGGCTGGGAGTGGTTGGAGAGGAGAATGAACCTACCTAAAGTC TT TAG AGG AAC ACG CGT CTG TAG TGG GCG TGG CCG TTT CAC TTC GGT TG

C

I S L C A D I T R T G K V K P

Genomic position	Amino acid change	Hetero/Homo	Disease	V-ID	Case/Control with this mutation	P-value	OR(95%CI)	Type of Study
Chr1 g.17662662T>C	NP_036519.2 p.L117		RA		733/735	0.00051	2.05 (1.51-2.86)	x2 test (Allelic)
Chr1 g.17662639T>C	NP_036519.2		RA		822/646	0.000008	1.97 (1.44-2.69)	x2 test (Allelic)
Chr1 g.17662639T>C	NP_036519.2		IDDM		1573/1732	0.87	1.01 (0.91-1.12)	Logistic regression (Allelic)

健常者の変異は黒、疾患に関わる変異は赤 (高さはP-value, もしくはstudy数に比例)



NGSの詳細表示

・NGS由来変異のクオリティコントロールのために、NGSの詳細を登録

Study name	Read depth	Mapping program	Detection of mutation	Quality level	Comment
study1	32	Bwa0.5.9	GATK	5	With indel realignment Default parameter
study2	20	SOAP2	SAMtool1.4	4	Without indel realignment Default Parameter
...					

Study name	mutation	Num of reads with mutations	Num of reads without mutations	Comment
study1	Chr 14 g.72684450G>A	25	30	With indel realignment Default parameter
study2	Chr14 g.72864480T>C	20	3	Without indel realignment Default parameter
...				

GWASデータ公開・共有方針

(統合DBプロジェクト疾患解析DB開発「倫理検討委員会」)

レベル1	頻度データ (遺伝子型、アレル、ハプロタイプ) SNPおよびCNV統計解析結果
データ アクセス	ウェブサイトにおいて閲覧可能 (公開データ) * 但しデータを大量取得する場合はレベル2と同様の申請を要する
レベル2	個体のCNVデータ、大量のレベル1データ
データ アクセス	氏名、職名、連絡先、使用目的、e-mailアドレス (原則、所属機関から発行されたアドレス) を記入して申請する
レベル3	個体の遺伝子型および生データ (共有データ) (適切な説明・同意が得られていることが前提)
データ アクセス	データアクセス申請書を提出し許可を受ける 使用期間に応じて、データ使用報告書を提出する

制限アクセスデータ(レベル3データ)*に関するデータ提供及びアクセスの手続き
 [*カテゴリーC GWAS遺伝子型データ(genotype data)およびカテゴリーD GWAS生データ(raw data)]

