



# KEGG: Kyoto Encyclopedia of Genes and Genomes

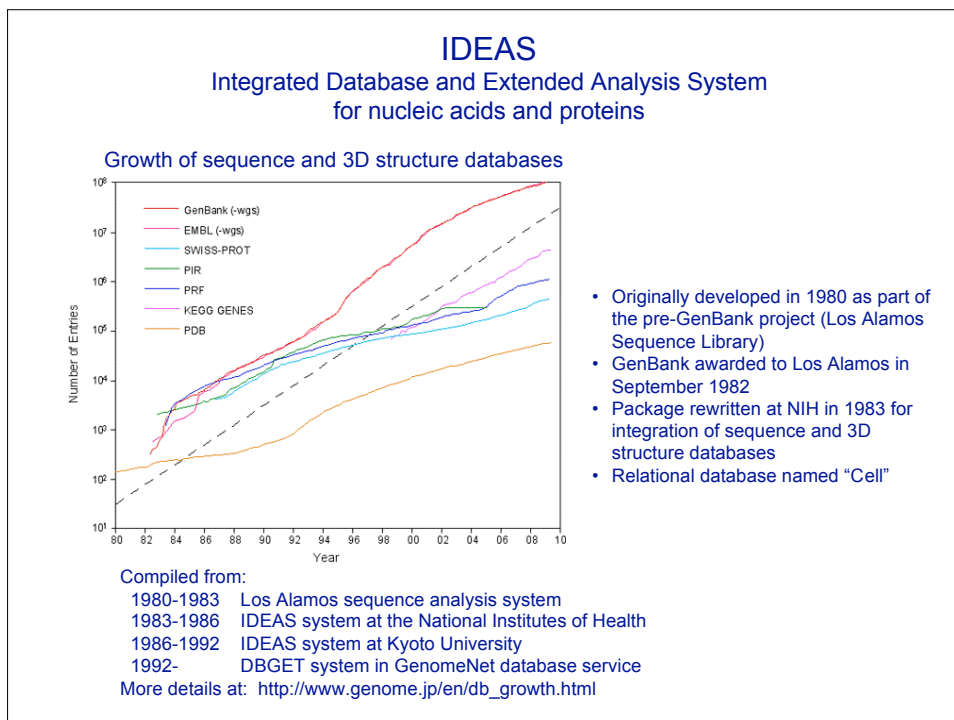
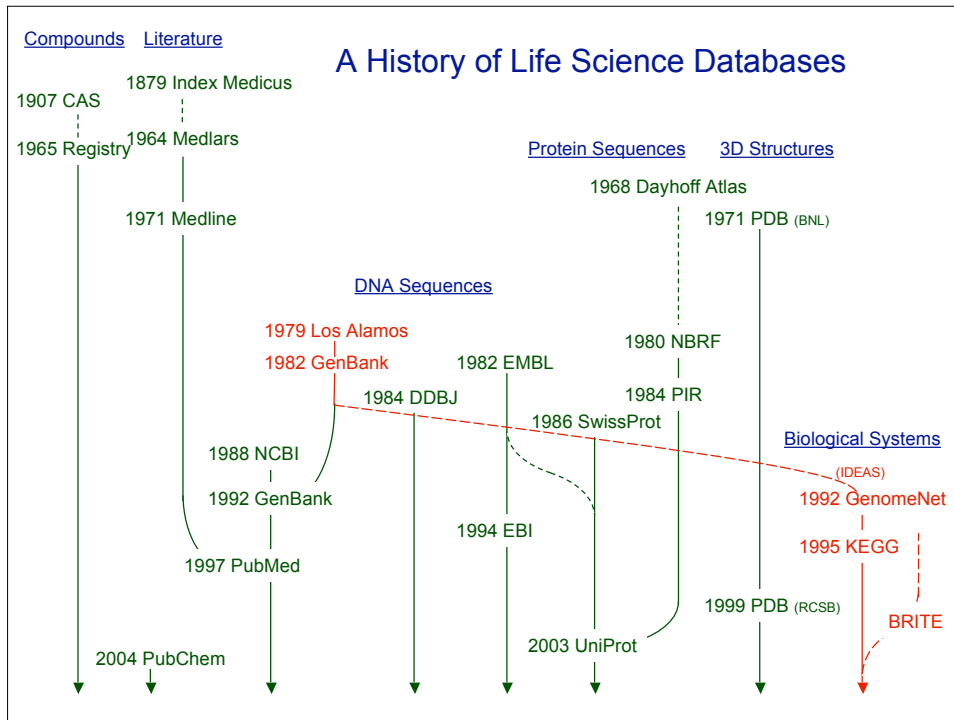
Minoru Kanehisa

*Bioinformatics Center, Institute for Chemical Research,  
Kyoto University  
Human Genome Center, Institute of Medical Science,  
University of Tokyo*

June 2009

## Background

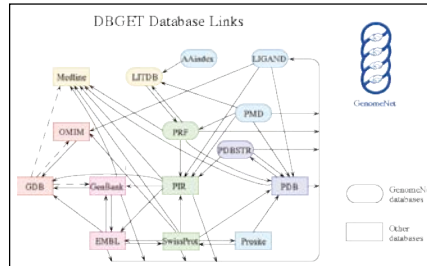
- Integration of life science databases
- Reconstruction of biological systems



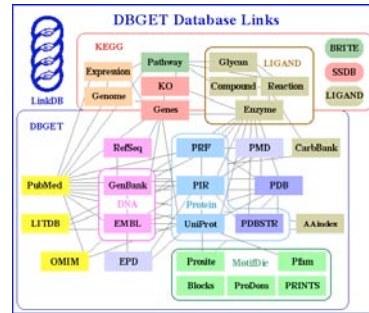
# DBGET

## Integrated Database Retrieval System for the Web of Molecular Biology Data

Link-based integration  
 Identifier database:entry  
 Relation database1:entry1 → database2:entry2



15 databases in September 1994



32 databases in September 2004

Category	bget	bfind	blink	brite	#DB
1. KEGG databases in DBGET	yes	yes	yes	yes	20
2. Other DBGET databases	yes	yes	yes	no	19
3. Searchable databases on the web	no	yes	yes	no	17
4. Link-only databases on the web	no	no	yes	no	94
5. PubMed database	yes	no	yes	no	1

151 databases in June 2009

# GenomeNet

<http://www.genome.jp/>

### Original Services

- KEGG
- DBGET
- MOTIF
- KAAS
- SIMCOMP / SUBCOMP
- KCaM

### Other Services

- BLAST / FASTA
- CLUSTALW / MAFFT / PRN
- CYORF / BSCORE

- Originally supported by Japanese Human Genome Project
- GenomeNet e-mail service launched in September 1992
- GenomeNet web service launched in July 1994
- Developed by Kanehisa Laboratories
- Operated by Kyoto University Bioinformatics Center
- Release history at: <http://www.genome.jp/en/release.html>

**KEGG**  
<http://www.genome.jp/kegg/>  
<http://www.kegg.jp/>

Organisms viewed as databases  
 Identifier organism:gene  
 Relation organism1:gene1 → organism2:gene2

- Collect only completely sequenced genomes for computer representation of cell/organism functions
- Organize higher-level biological and chemical knowledge to complement with the existing databases
- Depend on internal team members to develop high-quality databases (accept feedbacks, but no submissions)
- Develop associated research programs that utilize the new databases and demonstrate their utility
- Supported by multiple research grants to Minoru Kanehisa
- Project initiated in May 1995
- First release in December 1995

## What is Database?

	NCBI	KEGG
Database	Information infrastructure	Computer representation of biological systems
Content	Comprehensive repositories	Building blocks and wiring-diagrams
Integration	Linking	Reconstruction
Implementation	Entrez	KEGG
Main use	Individual data retrieval	Pathway mapping and Brite mapping

## Overview of KEGG

- From building blocks to biological systems
- Integration of genomics and chemistry

## Goal

An ultimate goal of **bioinformatics** is a complete computer representation of the cell and the organism, which will enable computational prediction of higher-level complexity, such as

- molecular interaction networks involving various cellular processes and
- phenotypes (morphological, physiological, and behavioral aspects) of entire organisms

from genomic information.

Kanehisa, M. and Bork, P.; Bioinformatics in the post-sequence era. *Nat. Genet.* **33**, 305-310 (2003).

## Building Blocks of Life Genomic and Chemical Spaces

### Genomic Space

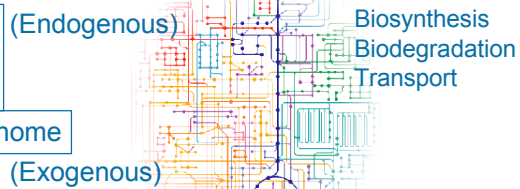
DNA (Gene)	Genome
RNA	Transcriptome
Protein	Proteome

Replication  
Transcription  
Translation



### Chemical Space

Metabolite	Metabolome
Glycan	Glycome
Lipid	Lipidome
Compound	Chemical genome



## Genomic and Chemical Spaces

### Genomic Space

- Contains genetic building blocks of life, i.e., DNA, RNA, and proteins
- Represents all possible sequences
- Uncovered by genomics, transcriptomics, and proteomics

### Chemical Space

- Contains chemical building blocks of life, e.g., small molecules, glycans, and lipids
- Represents all possible chemical structures
- Uncovered by metabolomics, glycomics, lipidomics, etc. for endogenous molecules
- Uncovered by chemical genomics for exogenous molecules

## Integration of Genomic and Chemical Spaces

### Genomic space

Repertoire of genes and proteins  
(genomics, proteomics, etc.)

Glycosyltransferase genes ↔  
Other enzyme gene groups ↔

Repertoire of genes in the biosphere  
(metagenomics)

Genetic factors of diseases ↔  
Drug targets ↔

Human ↔ Vector ↔ Pathogen

Human gut ↔ Microbial community

### Chemical space

Repertoire of endogenous molecules  
(metabolomics, glycomics, etc.)

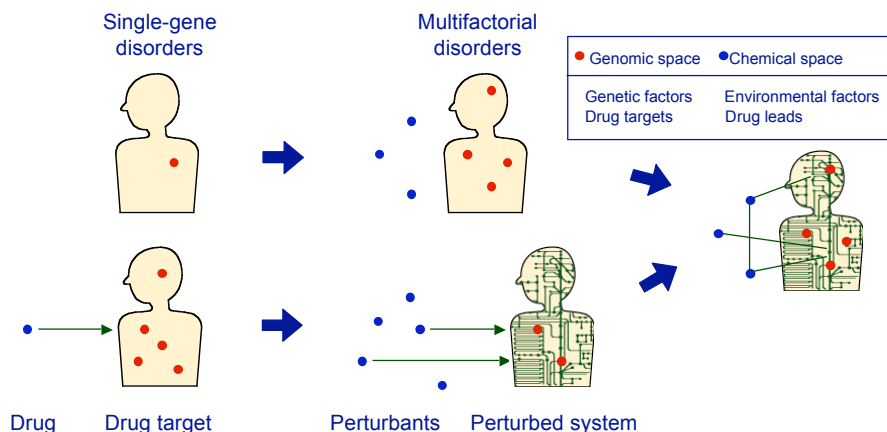
Glycan structures  
Lipids, secondary metabolites, etc.

Repertoire of exogenous molecules  
(chemical genomics)

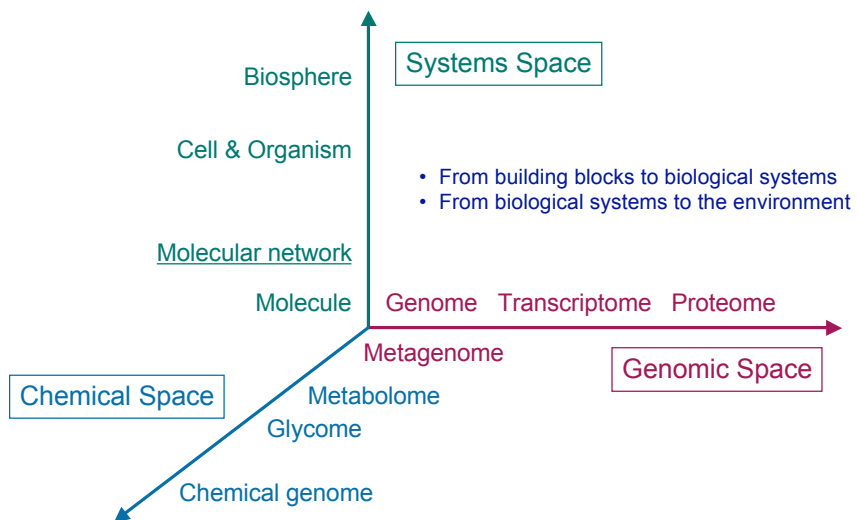
Environmental factors of diseases  
Drug leads

## Integration of Genomic and Chemical Spaces Medical and Pharmaceutical Implications

Diseases viewed as perturbed states of molecular systems  
Drugs viewed as perturbants to molecular systems

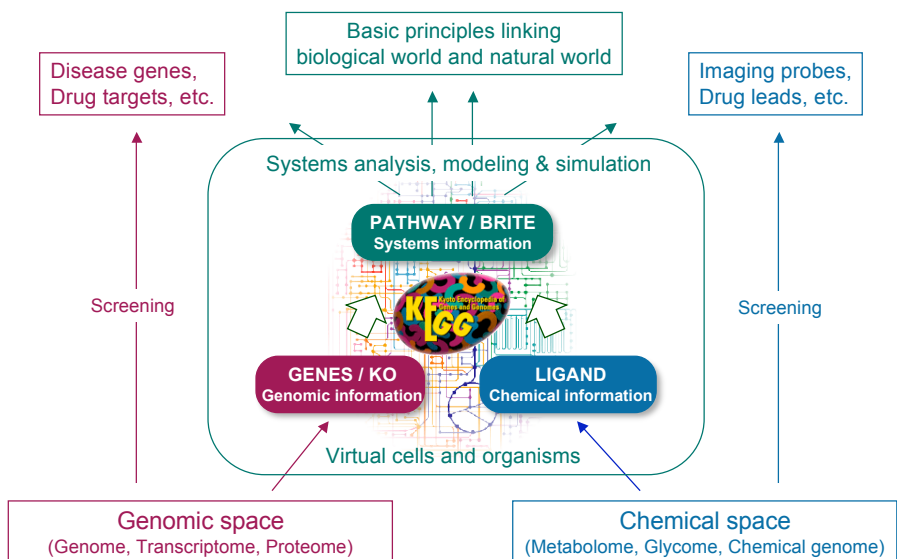


## Integration of Genomic and Chemical Spaces Bioinformatics Approaches to the Systems Space



Kanehisa, M. et al.; From genomics to chemical genomics: new developments in KEGG. *Nucl. Acids Res.* **34**, D354-D357 (2006).

## KEGG: Computer representation of biological systems



<http://www.genome.jp/kegg/>



## KEGG Databases

Database	Objects	Data size
KEGG PATHWAY	Pathway maps	330 (95,271)
KEGG BRITE	Functional hierarchies and ontologies	61 (21,684)
KEGG MODULE	Pathway modules and complexes	692
KEGG DISEASE	Diseases	104
KEGG ORTHOLOGY	KEGG Orthology (KO) groups	12,178
KEGG GENES	Genes in high-quality genomes	4,553,978
KEGG DGENES	Genes in draft genomes	151,662
KEGG EGENES	Genes as EST contigs	2,919,245
KEGG GENOME	Organisms	1,008 + 67
KEGG SSDB	Sequence similarities and best hit relations	-
KEGG COMPOUND	Metabolites and other small molecules	15,477
KEGG DRUG	Drugs	8,912
KEGG GLYCAN	Glycans	10,969
KEGG ENZYME	Enzyme nomenclature	5,061
KEGG REACTION	Biochemical reactions	7,915
KEGG RPAIR	Reactant pair chemical transformations	11,356
KEGG DPAIR	Drug pair chemical transformations	496

As of June 12, 2009

## KEGG Object Identifiers

Prefix + 5-digit number

Prefix	Content	Database	Example
K	Ortholog group	KEGG ORTHOLOGY	K04527 for insulin receptor
C	Chemical compound	KEGG COMPOUND	C00031 for D-glucose
D	Drug	KEGG DRUG	D01441 for Gleevec
G	Glycan	KEGG GLYCAN	G00109 for GM2
R	Reaction	KEGG REACTION	R00259 for EC 2.3.1.1
RP	Reactant pair	KEGG RPAIR	RP04458 for C00025_C00624
map/ko/ec/m/(org)	Pathway map	KEGG PATHWAY	hsa04930 for type II diabetes
br/jp/ko/(org)	Brite hierarchy	KEGG BRITE	ko01003 for glycosyltransferases
M	Pathway module	KEGG MODULE	M00008 for Entner-Doudoroff pathway
H	Human disease	KEGG DISEASE	H00004 for Chronic myeloid leukemia

Other identifiers

Identifier	Content	Database	Example
org	KEGG organism code <ul style="list-style-type: none"> <li>• Three letter code</li> <li>• Prefix d for draft genome</li> <li>• Prefix e for EST contig</li> </ul>	KEGG GENOME	hsa for Homo sapiens
org:gene	Gene entry	KEGG GENES	hsa:3643 for human insulin receptor
db:entry	Database entry	Any database	gn:hsa for Homo sapiens up:P06213 for UniProt insulin receptor



Try these examples by entering them in the search box of the KEGG top page.

## KEGG PATHWAY map for type II diabetes mellitus

**KEGG Pathway Map for Type II Diabetes Mellitus (Homo sapiens)**

**Gene Details (hsa:3643):**

Entry	3643	CDS	H. sapiens
Gene name	INSR		
Definition	Insulin receptor		
Orthology	K04527	Insulin receptor [EC:2.7.10.1]	
Pathway	PATH: hsa04520	Adherens junction	
	PATH: hsa04910	Insulin signaling pathway	
	PATH: hsa04930	Type II diabetes mellitus	
	PATH: hsa05050	Dentatorubropallidolysian atrophy (DRPLA)	
Class		(BRITE hierarchy)	
SSDB		(Orthology) (Paralogy) (Gene cluster)	
Motif	PFam: Recp_1_domain Purin-like TIL In3 Kinase_Tyr Kinase		
	PROSITE: PROSITE_KINASE_TYR_RECEPOR_TYR_KIN_II SHOTEN_KINASE_ATP_CYS_RICH_FK3 PROTEIN_KINASE_DOM		
Other DBs	OMIM: 147670		
	NCBI-GI: 119395734		
	NCBI-GeneID: 3643		
	NCIC: 4091		
	IRFDB: 40975		
	Ensembl: ENSG00000171105		
	UniProt: P06213		
LinkDB		(All DBs)	
Structure	PDB: 1B03 1P14 1IK3 3B04 2B05 1IKK 2B45 2B87 1X44 1RQQ 1GAG 2A08 228C 2D7G		
	Thumbnails		
Position	19p13.3-p13.2		
AA seq	1382 aa	(AA seq) (CD search)	
	ML003080AAALFLINVALLIQAAGELYPKVCVQMG1RNMLVLELEKNCV15QHL		
	ML003080AAALFLINVALLIQAAGELYPKVCVQMG1RNMLVLELEKNCV15QHL		
	ML003080AAALFLINVALLIQAAGELYPKVCVQMG1RNMLVLELEKNCV15QHL		
COMPOUND: C00031			
Entry	C00031	Compound	
Name	D-Glucose; Grape sugar; Dextrose		
Formula	C6H12O6		
Mass	180.0634		
Structure			
	<chem>C00031</chem>		
Remark	Same as: D00009		
Reaction	R00103 R00115 R00116 R00149 R00143 R00199 R00300 R00301 R00302 R00303 R00304 R00305 R00306 R00307 R00308 R00327		

## The KEGG Orthology (KO) System

KO (K number) entry for insulin receptor

**KEGG ORTHOLOGY: K04527**

Entry	K04527	KO
Name	INSR	
Definition	insulin receptor [EC:2.7.10.1]	
Class	Cellular Processes; Cell Communication; Adherens junction [PATH:ko04520] Cellular Processes; Endocrine System; Insulin signaling pathway [PATH:ko04910] Human Diseases; Neurodegenerative Diseases; [PATH:ko05050] Dentatorubropallidolysian atrophy (DRPLA) [PATH:ko05050] Human Diseases; Metabolic Disorders; Type II diabetes mellitus [PATH:ko04930] Protein Families; Metabolism; Protein Kinases [BR:ko01001] Protein Families; Cellular Processes and Signaling; Receptors and channels [BR:ko04000] Protein Families; Cellular Processes and Signaling; Cellular antigens [BR:ko04090]	
Other DBs	GO: 0005009	
Genes	BSA: 3643 (INSR) PTR: 455649 (INSR) MMU: 16337 (Insr) RNO: 24854 (Insr) CPA: 484990 (INSR) BTA: 408017 (INSR) MDO: 100027215 (LOC100027215) GGA: 420133 (INSR) DRE: 245699 (Insr) 245700 (Insr) SPU: 579152 (IR) NVE: NEMVE_v1q85808 (NEMVEDRAFT_v1q85808)	
LinkDB	(All DBs)	

KO identifiers (K numbers) represent manually defined ortholog groups corresponding to the KEGG pathway nodes and the BRITE hierarchy nodes (bottom leaves)

KEGG Orthology (ko00001)

- Metabolism
- Genetic Information Processing
- Environmental Information Processing
- Cellular Processes
- .....
- Cell Communication
- .....
- Endocrine System
- .....
- Human Diseases
- .....
- Neurodegenerative Diseases
- .....
- Metabolic Disorders
- .....

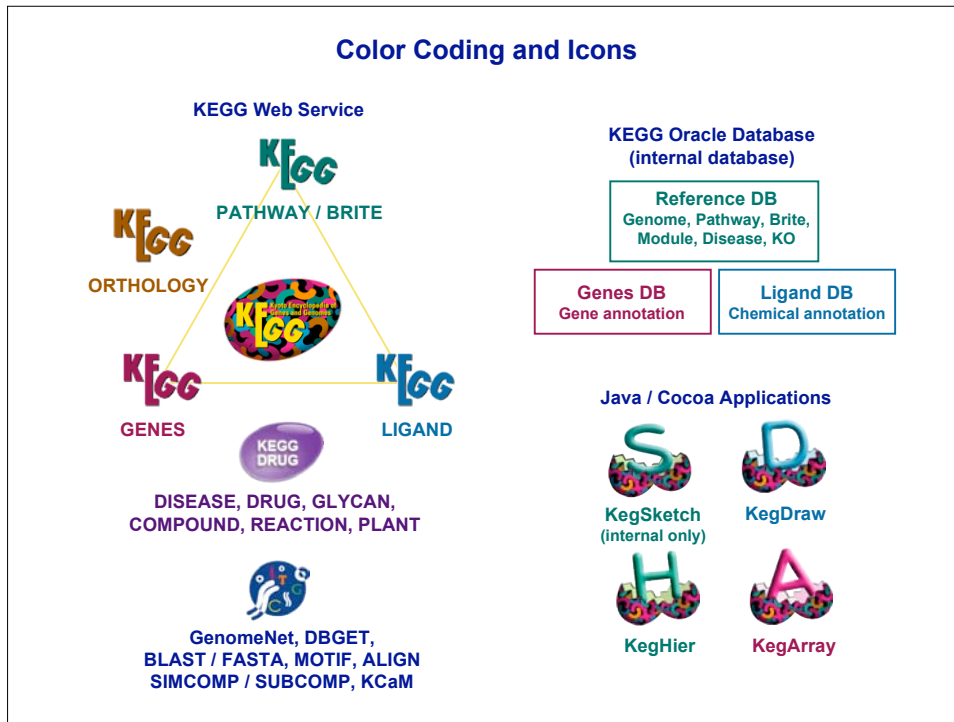
Protein Kinases (ko01001)

- Serine/threonine protein kinases: AGC group
- .....
- Tyrosine protein kinases
- .....
- InsR family
- .....
- Histidine protein kinases

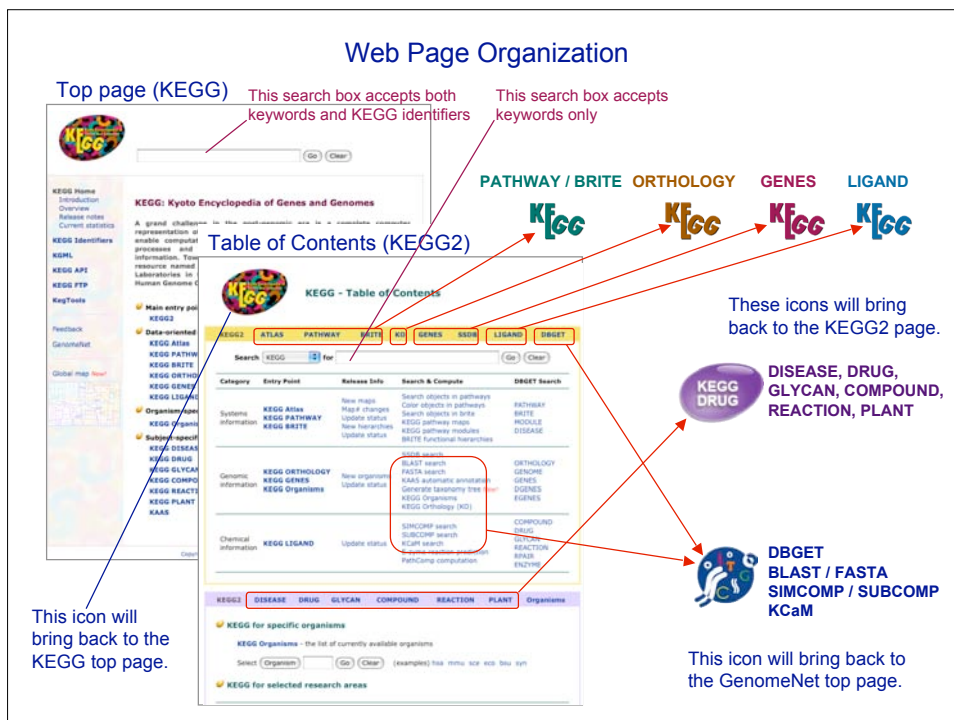
Receptors and Channels (ko04000)

- G-Protein Coupled Receptors
- 1 TM Signaling Receptors
- Receptor tyrosine kinase
- .....
- RTK class II
- .....
- Nuclear Receptors
- Ion Channels
- Other Channels
- .....
- Cellular antigens (ko04090)
- .....
- CD molecules
- Non-CD molecules

## Color Coding and Icons

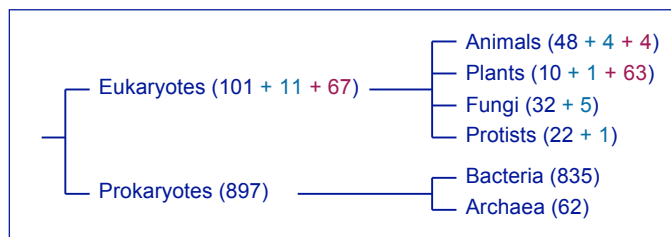


## Web Page Organization



# KEGG GENES and ortholog annotation

## KEGG Organisms



As of June 12, 2009

998 organisms in GENES - High-quality genomes with manual (koala) annotation

11 organisms in DGENES - Draft genomes with automatic (kaas) annotation

67 organisms in EGENES - EST contigs with automatic (kaas) annotation

Total number of genes/proteins in GENES 4,553,978

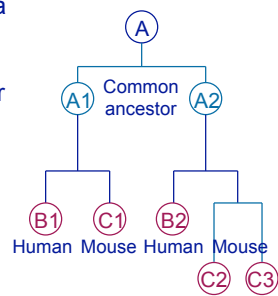
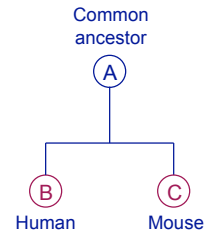
Total number of proteins in UniProt 8,385,695 (468,851 in SwissProt)

Fraction of UniProt covered by GENES about 90% (in terms of protein families)

KEGG GENES already covers most of the known protein universe

## Orthologs and Paralogs

- Sequence similarity between two genes (or proteins) may imply ortholog or paralog relationship.
- Orthologs are genes in different species evolved from a common ancestral gene by speciation and tend to have the same function.
- Paralogs are generated by gene duplication within a species and often represent diversified functions in a broader functional category.
- Identification of ortholog relationships is the basis for genome annotation (assigning gene functions), and it requires distinction from paralog relationships.

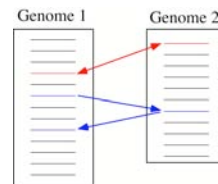


Orthologs: B–C, B1–C1  
 Co-orthologs: B2–(C2,C3)  
 Inparalogs: C2–C3  
 Outparalogs: B1–B2, B1–(C2,C3), B2–C1

## Computational Identification of Orthologs

### Between two species

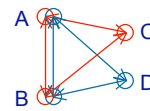
Bi-directional best hit (BBH)  
 (Reciprocal best hit)



### Among multiple species

#### 1. COG

Triangle of BBH relationships among three species



#### 2. KEGG OC

p-Quasi clique among multiple species

Superposition of ABC and ABD



Clique  
 (completely connected subgraph)

p-Quasi clique is an almost complete subgraph, where the degree of completeness is represented by p.

## Genome annotation in KEGG: KO (K number) assignment

### KEGG GENES

- Gene information for completely sequenced genomes
- Computationally generated from RefSeq and other public resources
- Manual annotation with KOALA and GFIT tools
- Automatic annotation by KAAS for draft genomes

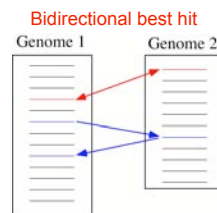
### KEGG GENES in Oracle

Genomes	997
Genes	4,553,977
Genes with KO	1,504,505
KO assignment	33%

As of June 12, 2009

### KEGG ORTHOLOGY (KO)

- Manually defined ortholog groups that correspond to KEGG pathway nodes and BRITE hierarchy nodes
- Identified by K numbers



### KEGG SSDB

- Sequence similarity scores and best hit relations
- Computationally generated from GENES by pairwise genome comparisons using SSEARCH

### KEGG OC

- Ortholog clusters
- Computationally generated from SSDB by a quasi-clique finding algorithm

## GFIT: Gene Function Identification Tool

DBHit(Length)	Over- lap (%)	Ident- (%)	Score	margin	Para- log	Cont.	KO	Orth	Current Annotation
> A mmu:20423(437)	434	68.0	1929	(295)	5		K06224	HH	[SP:SHH_MOUSE] Shh; sonic hedgehog ; K06224 hedgehog
> A gga:395615(425)	424	68.9	1921	(294)	6		K06224	HH	[SP:SHH_CHICK] SHH; sonic hedgehog homolog (Drosophila) ; K06224 hedgehog
> A rno:29499(437)	434	66.8	1892	(244)	7		K06224	HH	[SP:SHH_RAT] Shh; sonic hedgehog homolog (Drosophila) ; K06224 hedgehog
> A msc:716553(462)	462	65.6	1890	(259)	7		K06224	HH	LOC716553; sm K06224 hedgehog
> A hsa:6469(462)	462	64.7	1880	(252)	5		K06224	HH	[SP:A0247_HU] sonic hedgehog hedgehog
> A cfr:608860(461)	461	64.4	1878	(251)	6		K06224	HH	SHH, LOC608860; sonic hedgehog homolog (Drosophila) ; K06224 hedgehog
> A mdo:100016531(477)	477	62.7	1870	(194)	7		K06224	HH	LOC100016531; similar to potassium channel Kv beta 2.2 subunit; K06224 hedgehog
> A xla:398047(446)	438	63.5	1824	(178)	7		K06224	HH	[SP:SHH_XENLA] LOC398047; morphogen
> A ptr:743371(383)	418	64.4	1871	(46)	8		K06224	HH	SHH, LOC743371; sonic hedgehog homolog (Drosophila) ; K06224 hedgehog
> A xtr:100036603(396)	399	54.9	1433	(1330)	3		K06224	HH	[SP:ADHS7_XENTR] shh; desert hedgehog homolog
> A spu:373331(411)	418	48.3	1224	(1108)	17		K06224	HH	[SP:Q3L973_STRPU] hh; hedgehog ; K06224 hedgehog
> A dme:DMel_CG4637(471)	412	46.4	1182	(1079)	2		K06224	HH	[SP:A4V398_DROME] hh; hedgehog (EC:3.4.22.-); K06224 hedgehog
> A dpo:DPse_GA18321(472)	408	46.6	1167	(1056)	3		K06224	HH	[SP:Q29A9_DROME] GA18321 gene product from transcript GA18321-RA ; K06224 hedgehog
> A bta:286821(188)	187	89.8	1138	(557)	8		K06224	HH	SHH; sonic hedgehog homolog (Drosophila) ; K06224 hedgehog
> A nve:NMVE_v19241466(401)	360	45.3	1047	(287)	13		K06224	HH	[SP:A78Y14_NEMVE] NMVEDRAFT_v19241466; predicted protein; K06224 hedgehog
> A cel:TD5C12.10(1169)	193	30.6	279	(21)	14				[SP:Q21535_CAEEL] qua-1; QUAhog (hedgehog related)

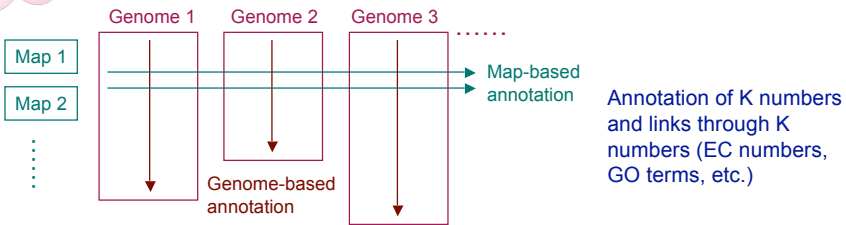
Color coding of BBH: red for the alignment with 80% or more overlap; yellow otherwise

- GFIT table is a summary of best hit relations for individual genes in a genome generated from the KEGG SSDB database. It has been used to manually assign K numbers (KO identifiers).
- KOALA processes this table and automatically assigns K numbers when considered safe to do so.





## KOALA: KEGG Orthology And Links Annotation



GRP	ORG	KEGG ID	KO (KOALA)	KO (GENES)	Orth (GENES)	OC	Score	Memo	Ann
E Ani	hsa	h564(448)	K04522	K04522	PSEN2	62906/6206/Animals.27398	18611		A
E Ani	ptr	ptr-457790(442)	K04522	K04522	PSEN2	62906/6206/Animals.27398	16345		A
E Ani	mcc	mcc-698770(454)	K04522	K04522	PSEN2	62906/6206/Animals.27398	10464		A
E Ani	mmu	mmu-19165(448)	K04522	K04522	PSEN2	62906/6206/Animals.27398	17031		A
E Ani	rno	rno-81751(448)	K04522	K04522	PSEN2	62906/6206/Animals.27398	16204		A
E Ani	cfa	cfa-490382(717)	K04522	K04522	PSEN2	62906/6206/Animals.27398	7851		A
E Ani	bta	bta-282010(449)	K04522	K04522	PSEN2	62906/6206/Animals.27398	16358		A
E Ani	ssc	ssc-780410(448)	K04522	K04522	PSEN2	62906/6206/Animals.27398	18558		A
E Ani	mdo	mdo-100026453(426)	K04522	K04522	PSEN2	62906/6206/Animals.27398	15132		A
E Ani	osa	osa-100076299(455)	K04522	>		62906/6206/Animals.27398	14136		A
E Ani	gga	gga-374188(451)	K04522	K04522	PSEN2	62906/6206/Animals.27398	10574		A
E Ani	xla	xla-397713(449)	K04522	K04522	PSEN2	62906/6206/Animals.27398	11327		A

KEGG ID	KO	Orth	Definition
h564	K04505/PSEN1	presenilin 1 (EC:3.4.23.-)	
h566	K04522/PSEN2	presenilin 2 (Alzheimer disease)	
ptr-457790	K04522/PSEN2	presenilin 2 (Alzheimer disease)	
mcc-698770	K04522/PSEN2	presenilin 2	
mmu-19165	K04522/PSEN1	presenilin 1 (EC:3.4.23.-)	
rno-81751	K04522/PSEN2	presenilin 2	
rno-29192	K04505/PSEN1	presenilin 1 (EC:3.4.23.-)	
cfa-490382	K04522/PSEN2	presenilin 2 (Alzheimer disease)	
bta-282010	K04522/PSEN2	presenilin 2 (Alzheimer disease)	
bta-282705	K04505/PSEN1	presenilin 1 (EC:3.4.23.-)	
ssc-780410	K04522/PSEN2	presenilin 2	
ssc-780411	K04505/PSEN1	presenilin 1 (EC:3.4.23.-)	
mdo-100019628	K04505/PSEN1	similar to presenilin 1	
mdo-100026453	K04522/PSEN2	similar to presenilin 2	
osa-100076299		similar to presenilin 2 (Alzheimer disease)	
osa-100091250		similar to presenilin 1	
gga-373977	K04505/PSEN1	presenilin 1 (Alzheimer disease)	
gga-374188	K04522/PSEN2	presenilin 2 (Alzheimer disease)	
xla-397713	K04522/PSEN2	presenilin-beta	
xla-399258	K04505/PSEN1	presenilin-alpha	

GFIT link      KOALA's Current suggestion assignment      OC links

## Public version of KEGG annotation tools

**KEGG ORTHOLOGY: K04522**

Entry: K04522 KO

Name: PSEN2, PS2

Definition: presenilin 2 [EC:3.4.23.-]

Class: Environmental Information Processing; Signal Transduction; Notch signaling pathway [PATH:ko08430]; Human Diseases; Neurodegenerative Diseases; Alzheimer's disease [PATH:ko08080]; Protein Families; Metabolism; Peptidases [BR:ko01002]

Genes: HSA: 5664(PSEN2), PTR: 457790(PSEN2), MCC: 698770(PSEN2), MGI: 19165(Psen2), RNO: 81751(Psen2), CFA: 490382(PSEN2), BTA: 282010(PSEN2), SSC: 780410(PSEN2), ECB: 100054506, MDO: 100026453, GGA: 374188(PSEN2), XLA: 397713, KTR: 549935(Psen2), DRE: 58026(Psen2), SPU: 578066

LinkDB: All DBs

**Read-only version of KOALA**

**KEGG Homo sapiens (human): 5664**

Entry: 5664 CDS h. sapiens

Gene name: PSEN2

Definition: presenilin 2 (Alzheimer disease 4)

Orthology: KO: K04522 presenilin 2 [EC:3.4.23.-]

Pathway: PATH: hsa04330 Notch signaling pathway; PATH: hsa08100 Alzheimer's disease

Class: Environmental Information Processing; Signal Transduction; Notch signaling pathway [PATH:hsa04330]; Human Diseases; Neurodegenerative Diseases; Alzheimer's disease [PATH:hsa08100]

SSDB: Orthology, Paralogs, Gene cluster, GFIT

Motif: Pfam: Presenilin VWA21 Herpes\_UL73 C4dC\_mol\_tron Peptidase\_A228 Dtd\_dored\_45\_h

Other DBs: OMIM: 600759, NCBI-GI: 156185679, NCBI-GeneID: 5664, HMC: 9589, HPI: 82868, Ensembl: ENSGN00000143801, UniProt: B1AP21\_P49810

LinkDB: All DBs

Position: 1a31-q42

AA seq: 448 aa

**Read-only version of GFIT**

**BLAST Search**

Enter query sequence: >hso:5664 PSEN2; presenilin 2 (Alzheimer disease 4); K04522 presenilin 2 [EC:3.4.23.-] (A)

Sequence data: >hso:5664 PSEN2; presenilin 2 (Alzheimer disease 4); K04522 presenilin 2 [EC:3.4.23.-] (A)

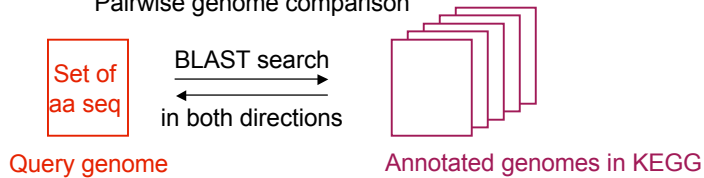
MLTMSADSEEEVCDERTLSMGEASPTPRSCQEC  
DRYVCSQVGRPPGLEELTLKYGARHVMFVY  
YTFYTEDTFSVGRLLKSLKLVLSKSYVIMFI  
FLFTYIYLVKLVNVMADYPTLLLVNMFAN  
LVFLKYLPERSAWVLLGASVYDVAVCPKQPL  
TVWAKLQSSVQALQLPDPQMEISYOSGEG  
KLGLGFYFVSLVQAAATSGDNNVTLACFVA  
TFGLIFYPSTDLVRLPMDLASHQLYI

Select program and database:

# KAAS: KEGG Automatic Annotation Server

<http://www.genome.jp/kegg/kaas/>

Pairwise genome comparison

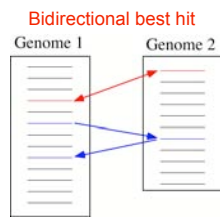


BLAST result screening by bi-directional best hit rate (BHR)

$$BHR = R_f \times R_r > 0.95$$

K number assignment by a heuristic scoring

$$S_{KO} = S_n - \log_2(mn) - \log_2\left(\sum_{k=N}^x C_k p^k (1-p)^{x-k}\right)$$



Moriya, Y., Itoh, M., Okuda, S., Yoshizawa, A.C., and Kanehisa, M.; KAAS: an automatic genome annotation and pathway reconstruction server. *Nucl. Acids Res.* 35, W182-W185 (2007).

## Genome Map Tools

KEGG Synechocystis sp. PCC6803

Genome page:

- Summary information
- Organism-specific pathway maps
- Organism-specific brite hierarchies
- Genome map browser

**Conserved gene cluster**

syn	001458	001454(K02067)	001453(K02049)	001451(K02050)	001450(K02051)	001456
syn	PCC7424_3464	PCC7424_3462	PCC7424_3524(K02049)	PCC7424_3526(K02050)	PCC7424_3527(K02051)	PCC7424_1680
mar	MAE_55270	MAE_53960	MAE_14770(K02049)	MAE_14790(K02050)	MAE_14800(K02051)	MAE_59360
cyt	Cyan7425_3581	Cyan7425_4366		Cyan7425_4571(K02050)	Cyan7425_4572(K02051)	Cyan7425_1598
cyp	PCC8801_4057	PCC8801_2463		PCC8801_4399(K02050)	PCC8801_4398(K02051)	PCC8801_2471
ava	Ava_4546	Ava_4544(K02067)		Ava_4542(K02049)	Ava_4541(K02050)	Ava_4540(K02051)
ava	001614	001611(K02067)	001611(K02049)	001609(K02050)	001606(K02051)	001605
cyb	CYB_02042	CYB_02040(K02067)		CYB_02037(K02049)	CYB_02036(K02050)	CYB_02035(K02051)
cya	CYA_0114	CYA_0619(K02067)		CYA_0615(K02049)	CYA_0614(K02050)	CYA_0613(K02051)
bel	Br1356	Br1355(K02067)		Br1352(K02049)	Br1351(K02050)	Br1350(K02051)
mig	Mig_1702(K02060)	Mig_1702(K02049)		Mig_1703(K02049)	Mig_1704	Mig_1665
cps	CPS_4945(K02072)			CPS_3314(K02049)	CPS_3315(K02050)	
trn	Trn_1199(K02072)			Trn_1185(K02049)	Trn_1184(K02050)	





## KEGG PATHWAY and BRITE: Reference knowledge base

### Data objects for computer representation of molecular systems

---

#### Element

gene, protein, small molecule, etc.

#### Pair (binary relation)

protein-protein interaction, drug-target relationship, etc.

#### Graph (wiring diagram)

pathway, complex, etc.

#### Simple list (membership)

pathway, complex, etc.

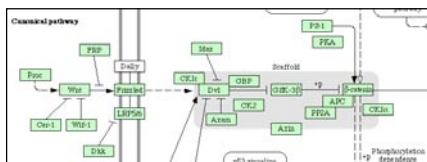
#### Hierarchical list

hierarchical classification, ontology, etc.

## Knowledge Representation of Systemic Functions

Graph (pathway map)

### KEGG PATHWAY



Simple list (membership)

### KEGG MODULE

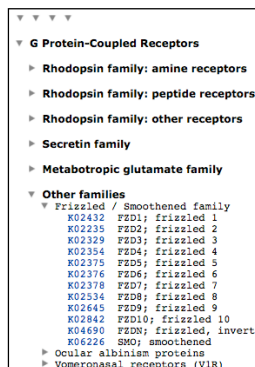
List of molecules that constitute a pathway, a complex, etc.

### KEGG DISEASE

List of disease genes, environmental factors, markers, drugs, etc.

Hierarchical list (ontology)

### KEGG BRITE



Data source: textbooks, review articles, other publications, specialists' websites

## KEGG PATHWAY Database

Collection of KEGG pathway maps

### Global Map

Metabolism Map (1)

### Metabolism

Carbohydrate Metabolism (16)

Energy Metabolism (8)

Lipid Metabolism (16)

Nucleotide Metabolism (2)

Amino Acid Metabolism (16)

Metabolism of Other Amino Acids (9)

Glycan Biosynthesis and Metabolism (15)

Biosynthesis of Polyketides and Nonribosomal Peptides (9)

Metabolism of Cofactors and Vitamins (12)

Biosynthesis of Secondary Metabolites (27)

Xenobiotics Biodegradation and Metabolism (26)

Overview (9)

### Genetic Information Processing

Transcription (2)

Translation (2)

Folding, Sorting and Degradation (8)

Replication and Repair (6)

### Environmental Information Processing

Membrane Transport (2)

Signal Transduction (14)

Signaling Molecules and Interaction (4)

### Cellular Processes

Cell Motility (3)

Cell Growth and Death (4)

Cell Communication (5)

Circulatory System (2)

Endocrine System (7)

Immune System (9)

Nervous System (2)

Sensory System (2)

Development (2)

Behavior (3)

### Human Diseases

Cancers (15)

Immune Disorders (6)

Neurodegenerative Diseases (4)

Metabolic Disorders (3)

Infectious Diseases (4)

### Drug Development

Chronology: Antibiotics (8)

Chronology: Antineoplastics (5)

Chronology: Nervous System Agents (9)

Chronology: Other Drugs (12)

Target Based Structure Classification (12)

Skeleton Based Structure Classification (8)

<http://www.genome.jp/kegg/pathway.html>

[http://www.genome.jp/kegg-bin/get\\_htext?br08901.keg](http://www.genome.jp/kegg-bin/get_htext?br08901.keg)

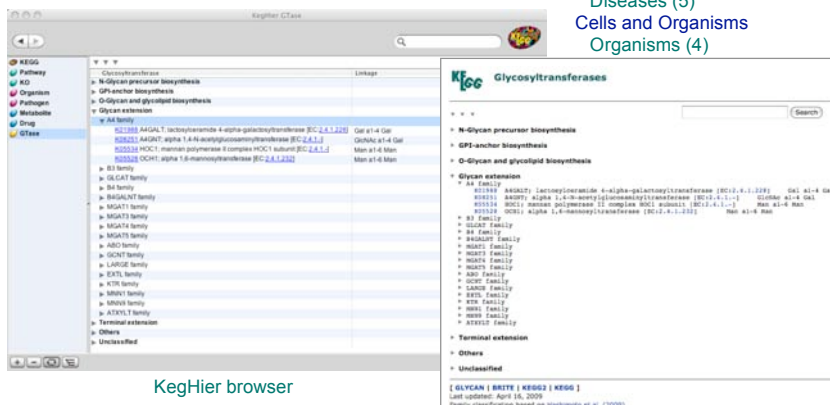
## KEGG BRITE Database Collection of BRITE functional hierarchies

### Genes and Proteins

- Network hierarchy (4)
- Protein families: metabolism (7)
- Protein families: genetic information processing (9)
- Protein families: environmental information processing (5)
- Protein families: cellular processes (7)

### Compounds and Reactions

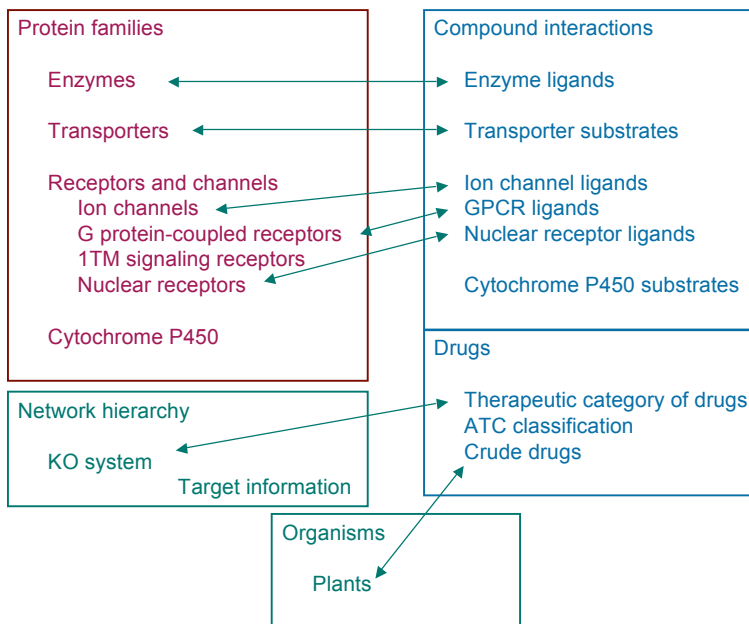
- Compounds (5)
  - Reactions (3)
  - Compound interactions (6)
- ### Drugs and Diseases
- Drugs (15)
  - Diseases (5)
- ### Cells and Organisms
- Organisms (4)



KegHier browser  
<http://www.genome.jp/kegg/brite.html>

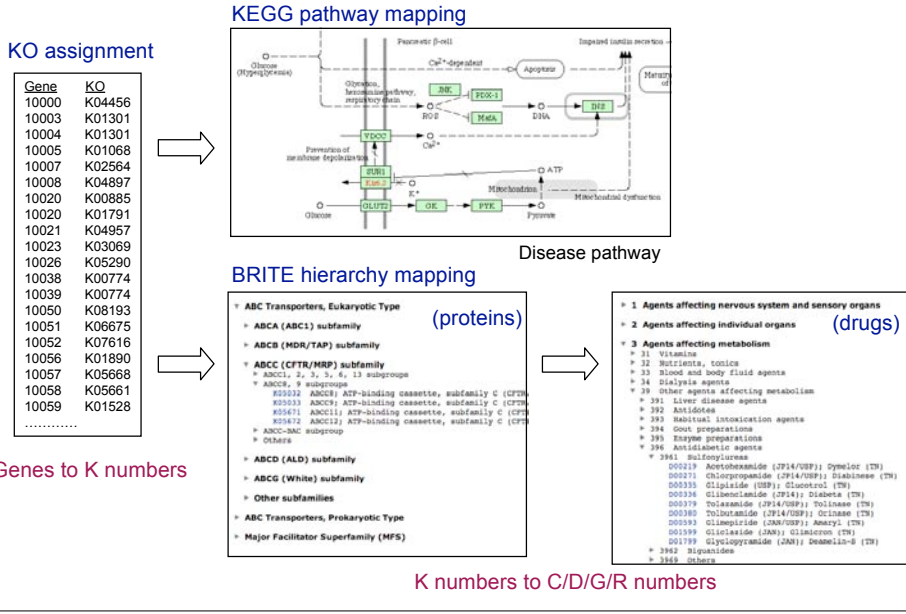
[http://www.genome.jp/kegg-bin/get\\_htext?br08902.keg](http://www.genome.jp/kegg-bin/get_htext?br08902.keg)

## Interaction/relation data linking BRITE hierarchies

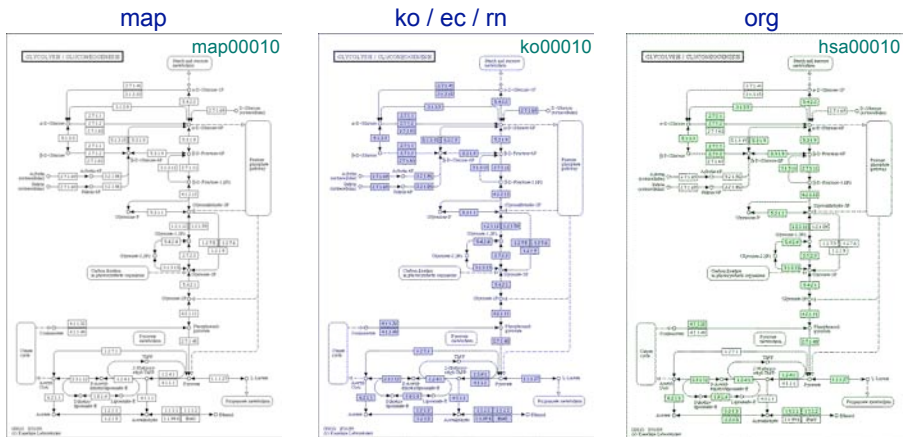


# Pathway mapping and BRITE mapping

## Linking genomes to biological systems and the environment



## Convention of map number prefix

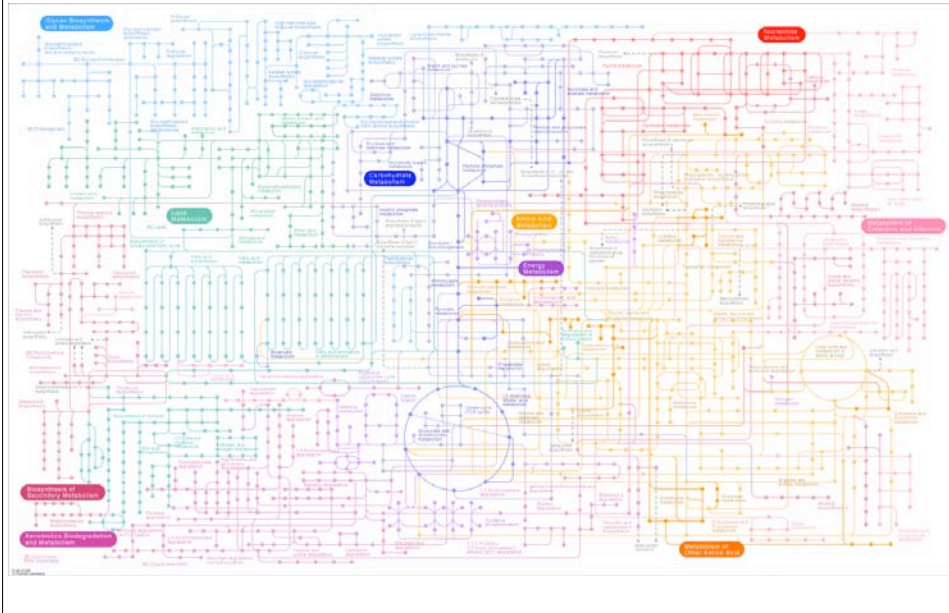


Each box is manually associated with KO identifier (K number), EC number, and reaction identifier (R number)

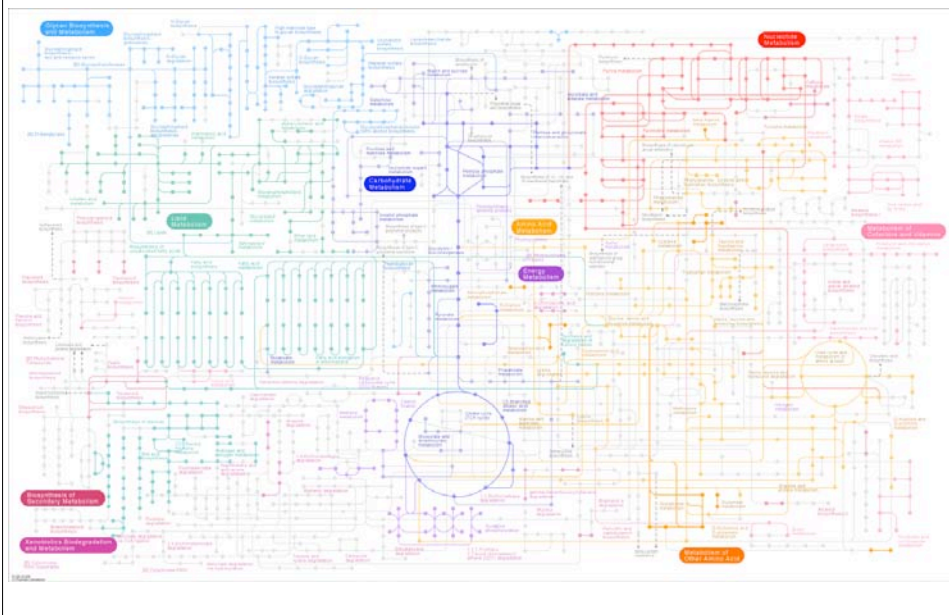
Blue boxes represent selection of K numbers, EC numbers, or R numbers

Green boxes correspond to gene identifiers in an organism that are computationally converted from K numbers

## Global metabolism map - Reference pathway

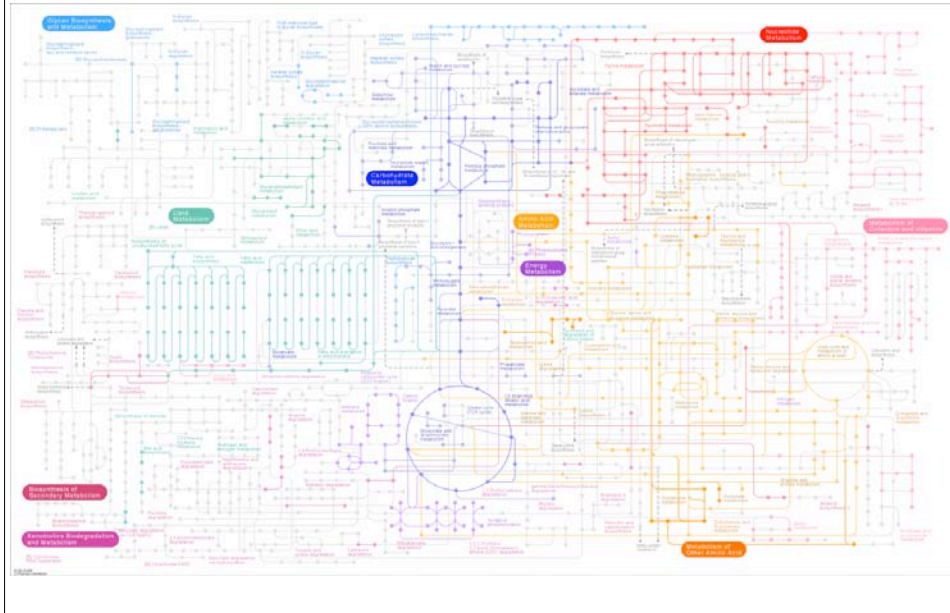


## Global metabolism map - Homo sapiens





## Global metabolism map - Escherichia coli



### Pathway Entry

Text-based representation of KEGG pathway maps

### Pathway MODULE

Tighter functional units in KEGG pathway maps

KEGG PATHWAY: ko00030

<b>Entry</b>	ko00030	Pathway																		
<b>Name</b>	Pentose phosphate pathway																			
<b>Description</b>	The pentose phosphate pathway is a process of glucose turnover that produces NADPH as reducing equivalents and pentoses as essential parts of nucleotides. There are two different phases in the pathway. One is irreversible oxidative phase in which glucose-6P is converted to ribulose-5P by oxidative decarboxylation, and NADPH is generated [MD:M00006]. The other is reversible non-oxidative phase in which phosphorylated sugars are interconverted to generate xylulose-5P, ribulose-5P, and ribose-5P [MD:M00007]. Phosphoribosyl pyrophosphate (PRPP) formed from ribose-5P [MD:M00005] is an activated compound used in the biosynthesis of histidine and purine/pyrimidine nucleotides. This pathway map also shows the Entner-Doudoroff pathway where 6-P-gluconate is dehydrated and then cleaved into pyruvate and glyceraldehyde-3P [MD:M00008].																			
<b>Class</b>	Metabolism; Carbohydrate Metabolism																			
<b>Pathway map</b>	ko00030 Pentose phosphate pathway																			
<b>Module</b>	<table border="1"> <tr><td>MD00004</td><td>Pen</td></tr> <tr><td>MD00005</td><td>Pen</td></tr> <tr><td>MD00006</td><td>Pen</td></tr> <tr><td>MD00007</td><td>Pen</td></tr> <tr><td>MD00008</td><td>Ent</td></tr> <tr><td>MD00009</td><td>Sen</td></tr> <tr><td>MD00010</td><td>P</td></tr> <tr><td>MD00011</td><td>Non</td></tr> <tr><td>MD00012</td><td>Pyru</td></tr> </table>		MD00004	Pen	MD00005	Pen	MD00006	Pen	MD00007	Pen	MD00008	Ent	MD00009	Sen	MD00010	P	MD00011	Non	MD00012	Pyru
MD00004	Pen																			
MD00005	Pen																			
MD00006	Pen																			
MD00007	Pen																			
MD00008	Ent																			
MD00009	Sen																			
MD00010	P																			
MD00011	Non																			
MD00012	Pyru																			
<b>Orthology</b>	<table border="1"> <tr><td>KR1810</td><td>glvA</td><td>glucose-6-phosphate isomerase [EC:5.3.1.9]</td></tr> <tr><td>KR1811</td><td>glvB</td><td>glucose-6-phosphate 1-dehydrogenase [EC:1.1.1.49]</td></tr> <tr><td>KR1812</td><td>glvC</td><td>6-phosphogluconolactonase [EC:3.1.1.31]</td></tr> <tr><td>KR1813</td><td>glvD</td><td>6-phosphogluconate dehydrogenase [EC:1.1.1.33]</td></tr> <tr><td>KR1814</td><td>glvE</td><td>6-phosphogluconate dehydrogenase [EC:1.1.1.44]</td></tr> <tr><td>KR1815</td><td>glvF</td><td>ribulose-phosphate 3-epimerase [EC:5.1.3.1]</td></tr> </table>		KR1810	glvA	glucose-6-phosphate isomerase [EC:5.3.1.9]	KR1811	glvB	glucose-6-phosphate 1-dehydrogenase [EC:1.1.1.49]	KR1812	glvC	6-phosphogluconolactonase [EC:3.1.1.31]	KR1813	glvD	6-phosphogluconate dehydrogenase [EC:1.1.1.33]	KR1814	glvE	6-phosphogluconate dehydrogenase [EC:1.1.1.44]	KR1815	glvF	ribulose-phosphate 3-epimerase [EC:5.1.3.1]
KR1810	glvA	glucose-6-phosphate isomerase [EC:5.3.1.9]																		
KR1811	glvB	glucose-6-phosphate 1-dehydrogenase [EC:1.1.1.49]																		
KR1812	glvC	6-phosphogluconolactonase [EC:3.1.1.31]																		
KR1813	glvD	6-phosphogluconate dehydrogenase [EC:1.1.1.33]																		
KR1814	glvE	6-phosphogluconate dehydrogenase [EC:1.1.1.44]																		
KR1815	glvF	ribulose-phosphate 3-epimerase [EC:5.1.3.1]																		

KEGG MODULE: M00008

<b>Entry</b>	M00008	Pathway	Module															
<b>Name</b>	Entner-Doudoroff pathway, glucose-6P → glyceraldehyde-3P + pyruvate																	
<b>Definition</b>	KR1810 K00030 K01857 K07484 K02698 K01625																	
<b>Pathway</b>	ko00030 Pentose phosphate pathway																	
<b>Orthology</b>	<table border="1"> <tr><td>KR1810</td><td>glucose-6-phosphate isomerase [EC:5.3.1.9]</td><td>[RN:R0739]</td></tr> <tr><td>K00030</td><td>glucose-6-phosphate 1-dehydrogenase [EC:1.1.1.49]</td><td>[RN:R0736]</td></tr> <tr><td>KR1812</td><td>6-phosphogluconolactonase [EC:3.1.1.31]</td><td>[RN:R07035]</td></tr> <tr><td>KR1813</td><td>6-phosphogluconate dehydrogenase [EC:1.1.1.33]</td><td>[RN:R07036]</td></tr> <tr><td>KR1625</td><td>2-dehydro-3-deoxyphosphogluconate aldolase [EC:4.1.1.2.14]</td><td>[RN:R05605]</td></tr> </table>			KR1810	glucose-6-phosphate isomerase [EC:5.3.1.9]	[RN:R0739]	K00030	glucose-6-phosphate 1-dehydrogenase [EC:1.1.1.49]	[RN:R0736]	KR1812	6-phosphogluconolactonase [EC:3.1.1.31]	[RN:R07035]	KR1813	6-phosphogluconate dehydrogenase [EC:1.1.1.33]	[RN:R07036]	KR1625	2-dehydro-3-deoxyphosphogluconate aldolase [EC:4.1.1.2.14]	[RN:R05605]
KR1810	glucose-6-phosphate isomerase [EC:5.3.1.9]	[RN:R0739]																
K00030	glucose-6-phosphate 1-dehydrogenase [EC:1.1.1.49]	[RN:R0736]																
KR1812	6-phosphogluconolactonase [EC:3.1.1.31]	[RN:R07035]																
KR1813	6-phosphogluconate dehydrogenase [EC:1.1.1.33]	[RN:R07036]																
KR1625	2-dehydro-3-deoxyphosphogluconate aldolase [EC:4.1.1.2.14]	[RN:R05605]																
<b>Reaction</b>	<table border="1"> <tr><td>R02798</td><td>C00048 → C01232</td></tr> <tr><td>R02799</td><td>C01232 → C01236</td></tr> <tr><td>R02800</td><td>C01236 → C00345</td></tr> <tr><td>R02801</td><td>C00345 → C04442</td></tr> <tr><td>R05605</td><td>C04442 → C00118 + C00022</td></tr> </table>			R02798	C00048 → C01232	R02799	C01232 → C01236	R02800	C01236 → C00345	R02801	C00345 → C04442	R05605	C04442 → C00118 + C00022					
R02798	C00048 → C01232																	
R02799	C01232 → C01236																	
R02800	C01236 → C00345																	
R02801	C00345 → C04442																	
R05605	C04442 → C00118 + C00022																	
<b>Module</b>	<table border="1"> <tr><td>C00048</td><td>alpha-D-Glucose 6-phosphate</td></tr> <tr><td>C01232</td><td>beta-D-Glucose 6-phosphate</td></tr> <tr><td>C01236</td><td>D-Glucono-1,5-lactone 6-phosphate</td></tr> <tr><td>C00345</td><td>6-Phospho-D-gluconate</td></tr> <tr><td>C04442</td><td>2-dehydro-3-deoxy-6-phospho-D-gluconate</td></tr> <tr><td>C00118</td><td>D-Glyceraldehyde 3-phosphate</td></tr> <tr><td>C00022</td><td>Pyruvate</td></tr> </table>			C00048	alpha-D-Glucose 6-phosphate	C01232	beta-D-Glucose 6-phosphate	C01236	D-Glucono-1,5-lactone 6-phosphate	C00345	6-Phospho-D-gluconate	C04442	2-dehydro-3-deoxy-6-phospho-D-gluconate	C00118	D-Glyceraldehyde 3-phosphate	C00022	Pyruvate	
C00048	alpha-D-Glucose 6-phosphate																	
C01232	beta-D-Glucose 6-phosphate																	
C01236	D-Glucono-1,5-lactone 6-phosphate																	
C00345	6-Phospho-D-gluconate																	
C04442	2-dehydro-3-deoxy-6-phospho-D-gluconate																	
C00118	D-Glyceraldehyde 3-phosphate																	
C00022	Pyruvate																	

### KEGG-NCBI Collaboration

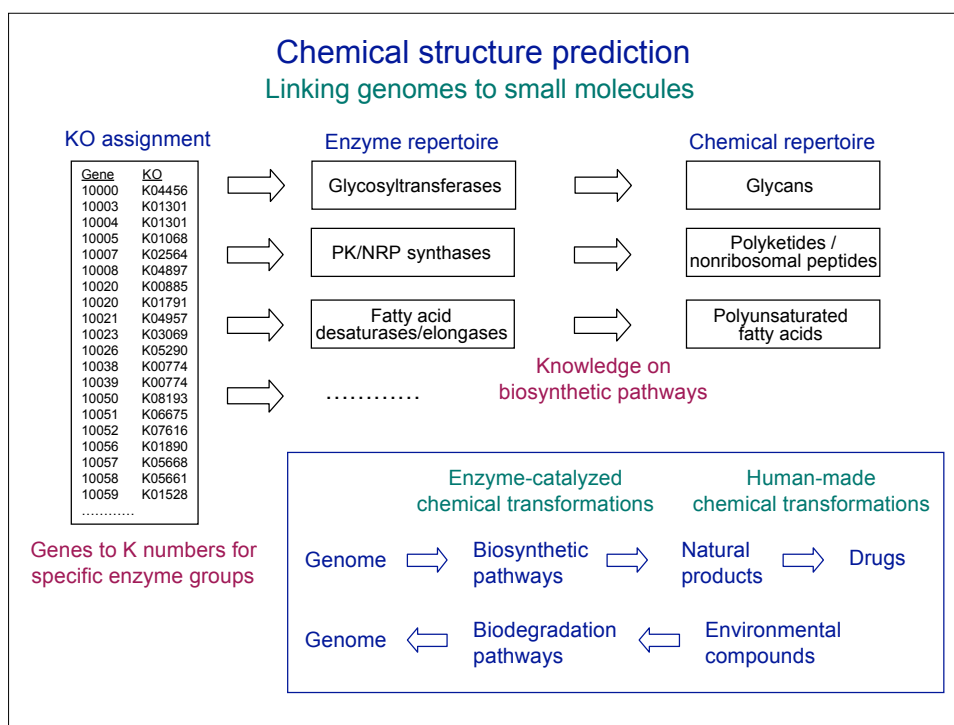
KEGG	NCBI
GENES	RefSeq (Gene)
PATHWAY	BioSystems
COMPOUND	PubChem
DRUG	PubChem



## KEGG LIGAND and chemical bioinformatics

### Bioinformatics for Small Molecules

1. Chemical structure similarity
  - Comparison of bit-represented vectors (fingerprints)
  - Comparison of graph objects
2. Chemical building blocks
  - Conserved substructures as building blocks of compounds/drugs
  - Variable substructures as building blocks of reactivity/efficacy
3. Network modules
  - Genomic module, e.g. operon
  - Chemical module, e.g. overall reaction
4. Predictive methods
  - Interaction prediction, e.g. toxicity
  - Reaction prediction, e.g. metabolic fate
5. Examples
  - Plant/fungi/bacterial genomes and secondary metabolites via biosynthetic pathways
  - Bacterial genomes and environmental compounds via biodegradation pathways



### Linking genomes to chemical structures

**Glycans**

- Hashimoto, K., Goto, S., Kawano, S., Aoki-Kinoshita, K.F., Ueda, N., Hamajima, M., Kawasaki, T., and Kanehisa, M.; KEGG as a glycome informatics resource. *Glycobiology* 16, 63R-70R (2006).
- Kawano, S., Hashimoto, K., Miyama, T., Goto, S., and Kanehisa, M.; Prediction of glycan structures from gene expression data based on glycosyltransferase reactions. *Bioinformatics* 21, 3976-3982 (2005).
- Hashimoto, K., Tokimatsu, T., Kawano, S., Yoshizawa, A.C., Okuda, S., Goto, S., and Kanehisa, M.; Comprehensive analysis of glycosyltransferases in eukaryotic genomes for structural and functional characterization of glycans. *Carbohydrate Res.* 344, 881-887 (2009).

**Polyketides & nonribosomal peptides**

- Minowa, Y., Araki, M., and Kanehisa, M.; Comprehensive analysis of distinctive polyketide and nonribosomal peptide structural motifs encoded in microbial genomes. *J. Mol. Biol.* 368, 1500-1517 (2007).

**Polyunsaturated fatty acids**

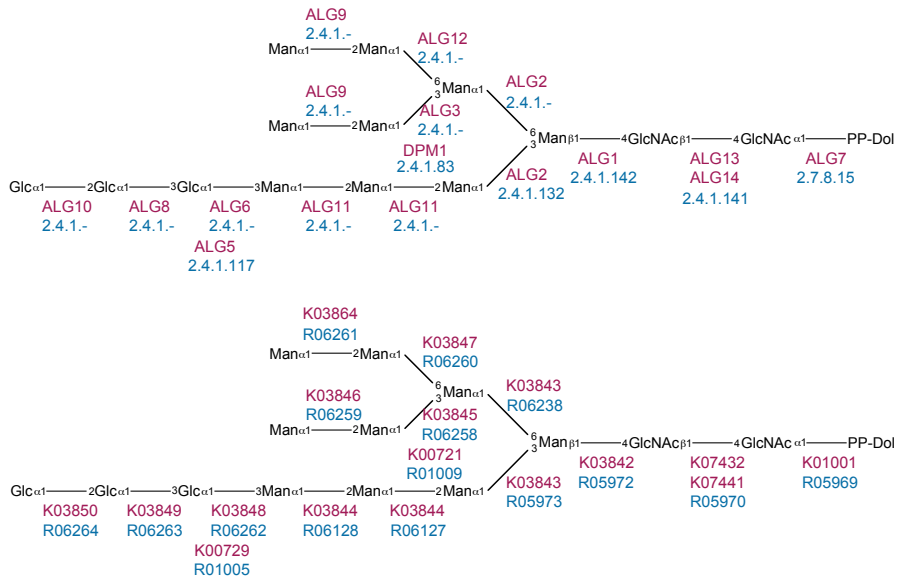
- Hashimoto, K., Yoshizawa, A.C., Okuda, S., Kuma, K., Goto, S., and Kanehisa, M.; The repertoire of desaturases and elongases reveals fatty acid variations in 56 eukaryotic genomes. *J. Lipid Res.* 49, 183-191 (2008).

**Environmental compounds**

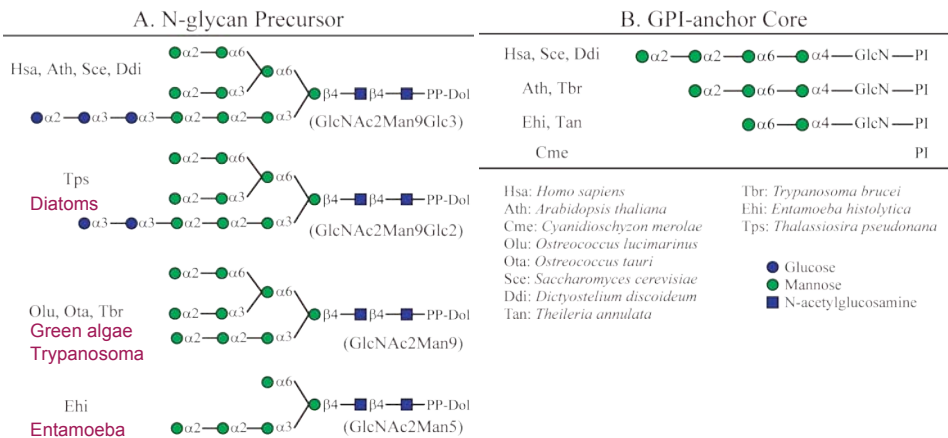
- Oh, M., Yamada, T., Hattori, M., Goto, S., and Kanehisa, M.; Systematic analysis of enzyme-catalyzed reaction patterns and prediction of microbial biodegradation pathways. *J. Chem. Info. Model.* 47, 1702-1712 (2007).



## N-Glycan Precursor Structure

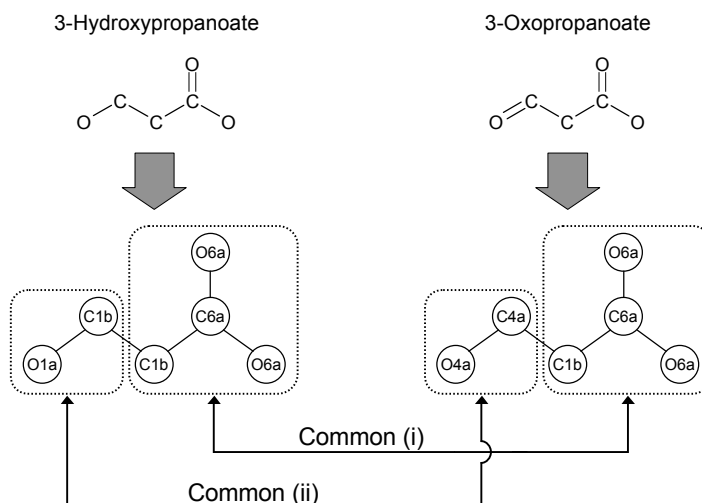


## Truncated glycan structures in parasitic protists and algae predicted from genomic information



Hashimoto, K., Tokimatsu, T., Kawano, S., Yoshizawa, A.C., Okuda, S., Goto, S., and Kanehisa, M.; Comprehensive analysis of glycosyltransferases in eukaryotic genomes for structural and functional characterization of glycans. *Carbohydrate Res.* 344, 881-887 (2009).

## Chemical structure comparison based on atom typing



Hattori, M., Okuno, Y., Goto, S., and Kanehisa, M.; Development of a chemical structure comparison method for integrated analysis of chemical and genomic information in the metabolic pathways. *J. Am. Chem. Soc.* 125, 11853-11865 (2003).

## KEGG atom types

### Carbon 23 types

Alkane	C1a	R-CH3
	C1b	R-CH2-R
	C1c	R-CH(R)-R
	C1d	R-C(R)2-R
Cyclic alkane	C1x	ring-CH2-ring
	C1y	ring-CH(R)-ring
	C1z	ring-CH(R)2-ring
Alkene	C2a	R=CH2
	C2b	R=CH-R
	C2c	R=C(R)2
Cyclic alkene	C2x	ring-CH=ring
	C2y	ring-C(R)=ring ring-C(=R)-ring
Alkyne	C3a	R=CH
	C3b	R=C-R
Aldehyde	C4a	R-CH=O
Ketone	C5a	R-C(=O)-R
Cyclic ketone	C5x	ring-C(=O)-ring
Carboxylic acid	C6a	R-C(=O)-OH
Carboxylic ester	C7a	R-C(=O)-O-R
	C7x	ring-C(=O)-O-ring
Aromatic ring	C8x	ring-CH=ring
	C8y	ring-C(R)=ring
Undefined C	C0	

### Nitrogen 16 types

Amine	N1a	R-NH2
	N1b	R-NH-R
	N1c	R-N(R)2
	N1d	R-N(R)3+
Cyclic amine	N1x	ring-NH-ring
	N1y	ring-N(R)-ring
Imine	N2a	R=N-H
	N2b	R=N-R
Cyclic imine	N2x	ring-N=ring
	N2y	ring-N(R)=ring
Cyan	N3a	R#N
Aromatic ring	N4x	ring-NH-ring
	N4y	ring-N(R)-ring
	N5x	ring-N=ring
	N5y	ring-N(R)=ring
Undefined N	N0	

### Oxygen 18 types

Hydroxy	O1a	R-OH
	O1b	N-OH
	O1c	P-OH
	O1d	S-OH
Ether	O2a	R-O-R
	O2b	P-O-R
	O2c	P-O-P
	O2x	ring-O-ring
Oxo	O3a	N=O
	O3b	P=O
	O3c	S=O
Aldehyde	O4a	R-CH=O
Ketone	O5a	R-C(=O)-R
	O5x	ring-C(=O)-ring
Carboxylic acid	O6a	R-C(=O)-OH
Ester	O7a	R-C(=O)-O-R
	O7x	ring-C(=O)-O-ring
Undefined O	O0	

### Sulfur 7 types

Thiol	S1a	R-SH
Thioether	S2a	R-S-R
	S2x	ring-S-ring
Disulfide	S3a	R-S-S-R
	S3x	ring-S-S-ring
Sulfate	S4a	R-SO3
Undefined S	S0	

### Phosphorus 2 types

Attached to other elements	P1a	P-R
Attached to oxygen	P1b	P-O

### Other elements 2 types

Halogens	X	F, Cl, Br, I
Others	Z	

## Reactant pairs extracted from enzymatic reactions

### EC1. Oxidoreductases



### EC2. Transferases



### EC3. Hydrolases



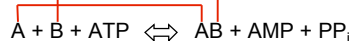
### EC4. Lyases



### EC5. Isomerases



### EC6. Ligases



### EC (Enzyme Commission) number

- EC numbers represent classification of enzymatic reactions (and enzymes)
- Manually assigned by IUBMB-IUPAC Biochemical Nomenclature Committee
- Published experimental evidence on the enzyme and the reaction is required
- Four numbers separated by periods, representing reaction specificity (class, sub-class, and sub-subclass numbers) and substrate specificity (serial number).

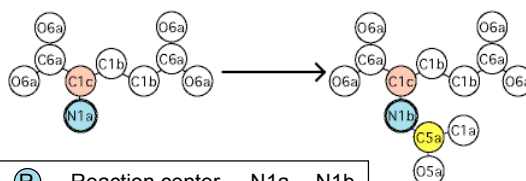
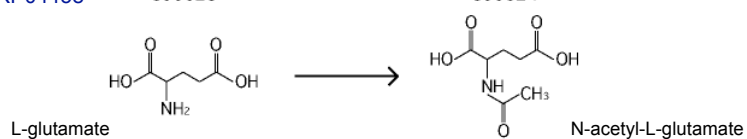
## Reactant pair transformation patterns (RDM patterns)

R00259 C00024 + C00025  $\rightleftharpoons$  C00010 + C00624 [EC:2.3.1.1]  
Acetyl-CoA + L-Glutamate  $\rightleftharpoons$  CoA + N-Acetyl-L-glutamate

RP04458

C00025

C00624



<span style="border: 1px solid blue; border-radius: 50%; padding: 2px;">R</span>	Reaction center	N1a -- N1b
<span style="border: 1px solid yellow; border-radius: 50%; padding: 2px;">D</span>	Difference atom	(H) -- C5a
<span style="border: 1px solid red; border-radius: 50%; padding: 2px;">M</span>	Matched atom	C1c -- C1c

RDM pattern

Kotera, M., Okuno, Y., Hattori, M., Goto, S., and Kanehisa, M.; Computational assignment of the EC numbers for genomic-scale analysis of enzymatic reactions. *J. Am. Chem. Soc.* 126, 16487-16498 (2004).

## Prediction of xenobiotics degradation pathways

Predicted pathways for 1,2,3,4-tetrachlorobenzene

Carbohydrate  
Energy  
Lipid  
Nucleotide  
Amino Acid  
Other aa  
Glycan  
PK/NRP  
Cofactors  
Sec.metab.  
Xenobiotics

gamma-tetachlorocyclohexane pathway (map00861)

RDM patterns for xenobiotics biodegradation pathways

RDM patterns  
Chemical structure transformation patterns in known enzymatic reactions

Oh, M., Yamada, T., Hattori, M., Goto, S., and Kanehisa, M.; Systematic analysis of enzyme-catalyzed reaction patterns and prediction of microbial biodegradation pathways. *J. Chem. Info. Model.* 47, 1702-1712 (2007).

## Structural Similarity Search

**KcAM** COMPOUND: C04884

**Entry:** C04884  
**Name:** 1,4-bis(2-acetyl-3-O-galactosylamino)-3-O-galactosyl-β-D-glucosylceramide  
**Formula:** C58H98N10O34

**Reaction:** C04884 → C04885  
**Enzyme:** EC 3.2.1.152

**Similarity Search Results:**

No.	Entry	Structure	Similarity Score
1	C04884		1.00
2	C04885		0.99
3	C04886		0.98
4	C04887		0.97

**KcAM** GLYCAN: G00109

**Entry:** G00109  
**Name:** Glycan  
**Structure:** GalNAc6S6 → GalNAc6S6 → GlcNAc6 → Cer  
**Class:** Glycolipid, sphingolipid

**Similarity Search Results:**

No.	Entry	Structure	Similarity Score
1	G00109		1.00
2	G00110		0.99
3	G00111		0.98
4	G00112		0.97
5	G00113		0.96

SIMCOMP for chemical structures of small molecules

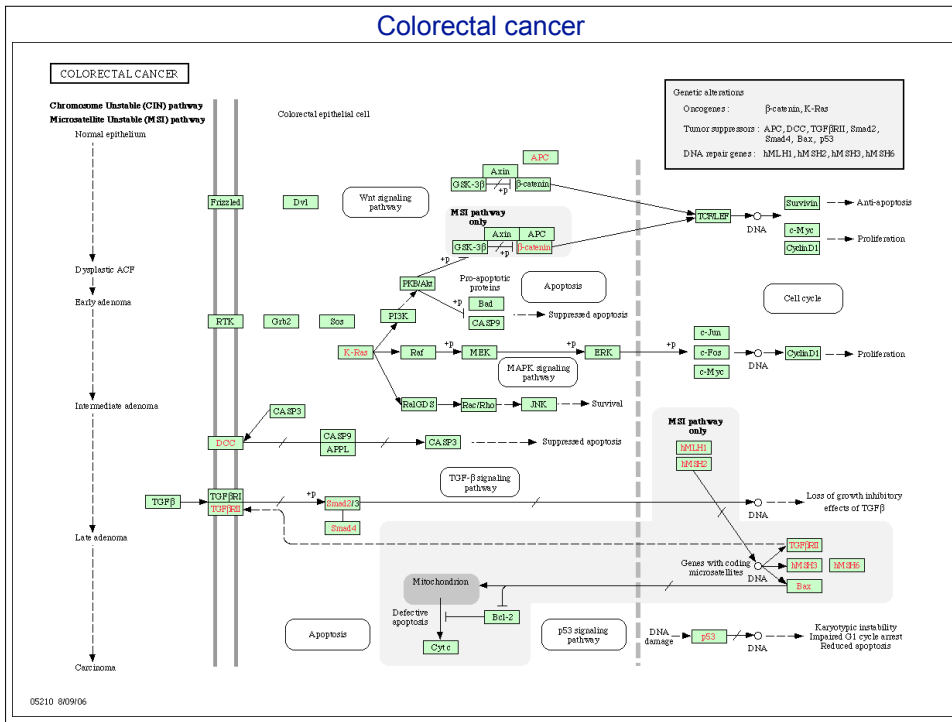
KcAM for glycan structures

## KEGG DISEASE and DRUG: An ultimate knowledge base

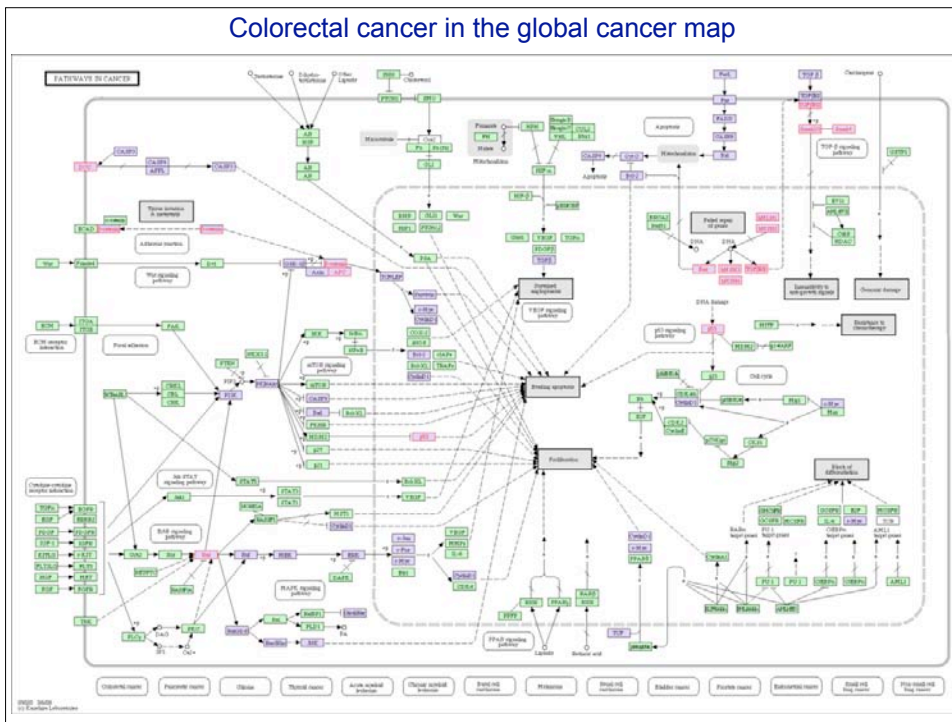
### Disease and Drug Information in KEGG

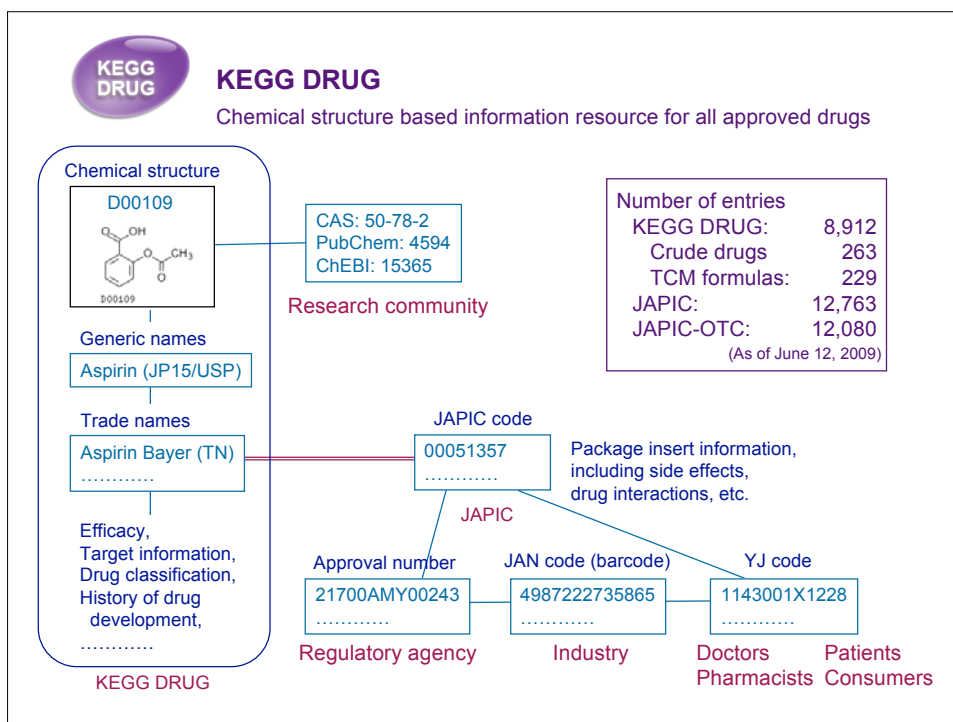
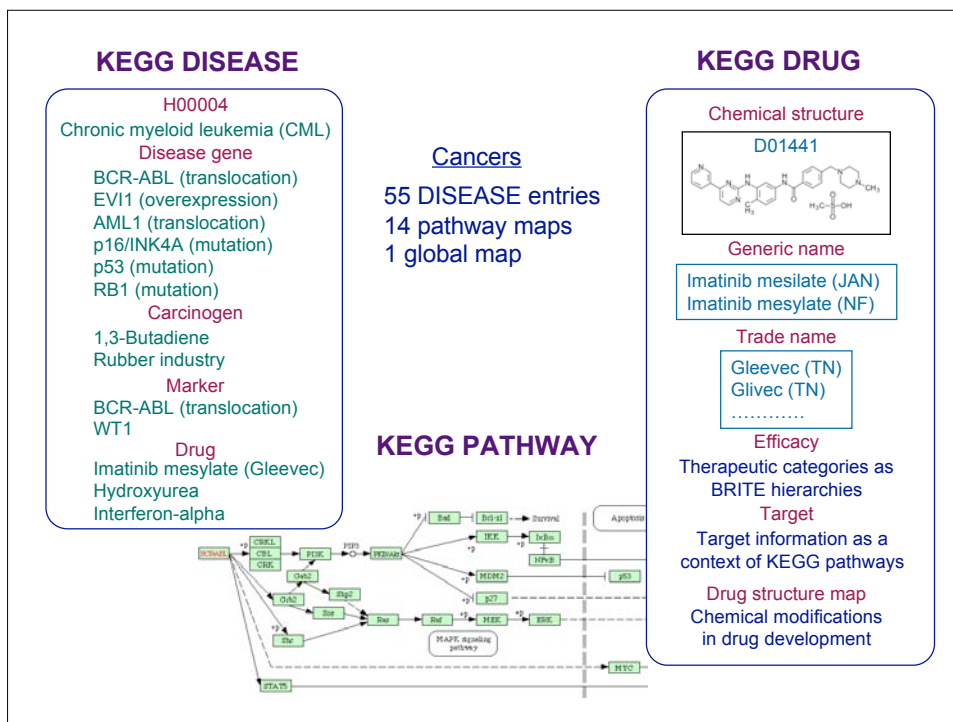
	KEGG DISEASE	KEGG DRUG
URL	<a href="http://www.genome.jp/kegg/disease/">http://www.genome.jp/kegg/disease/</a>	<a href="http://www.genome.jp/kegg/drug/">http://www.genome.jp/kegg/drug/</a>
Content	Lists of disease genes and molecular factors	Chemical structure based collection of all approved drugs in Japan, USA, and Europe
Pathway	KEGG pathway maps for human diseases: cancers, immune disorders, neurodegenerative diseases, metabolic disorders, and infectious diseases	KEGG DRUG structure maps for drug development
BRITE hierarchy	Disease classifications including: Pathogens and infectious diseases Human diseases ICD-10 disease classification	Drug classifications including: Therapeutic category of drugs (Japan) USP drug classification (USA) ATC classification (WHO) TCM (Traditional Chinese Medicine) drugs in Japan Crude drugs

## Colorectal cancer



## Colorectal cancer in the global cancer map







# Herbal Medicine in KEGG

**KEGG Traditional Chinese Medicine in Japan**

Crude drug: combination of compounds  
**Kakkon (Pueraria root)**  
 DRUG: D06693

- ▼ Crude Drugs
  - ▶ Drugs for relieving the exterior (diaphoretics)
    - ▶ Diaphoretics pungent in flavor and warm in property
      - ▶ **Diaphoretics pungent in flavor and cool in property**
        - D06693 Sweetroot root (JP15); Pueraria root (TN)
        - D06744 Chrysanthemum flower (JP15)
        - D06723 Burdock fruit (JP15)
        - D06727 Suppleurum root (JP15); Suppleurum root (TN)
        - D06745 Cimicifuga rhizome (JP15); Cimicifuga rhizome
        - D05431 Peppermint (JP15/NF); Peppermint (TN)
  - ▶ Drugs for clearing heat
  - ▶ Purgative drugs
  - ▶ Drugs for dispelling wind-damp (antirheumatics)
  - ▶ Aromatic drugs for resolving dampness
  - ▶ Drugs for inducing diuresis and excreting dampness (diuretics)
  - ▶ Drugs for warming the interior
  - ▶ Drugs for regulating Qi
  - ▶ Drugs for removing food stagnation (digestants)
  - ▶ Drugs for expelling parasites (anthelmintics)
  - ▶ Hemostatic drugs
  - ▶ Drugs for promoting blood circulation and removing blood stasis
  - ▶ Expectorant, antitussive and anti-asthmatic drugs
  - ▶ Drugs for calming the spirit (tranquillizers)
  - ▶ Drugs for calming the liver and extinguishing wind
  - ▶ Drugs for resuscitation
  - ▶ Drugs for supplementing efficiency (tonics)
  - ▶ Astringent drugs
  - ▶ Emetic drugs
  - ▶ Drugs for external use

Entry	D06693	Crude Drug
<b>Name</b>	Pueraria root (JP15); Pueraria root (TN)	<b>Biosynthetic pathways</b>
<b>Component</b>	Starch [DR:D06507] [CPD:C00369], Daidzin [CPD:C10216], Daidzein [CPD:C10208], Puerarin [CPD:C10524], Puerarinxyloside, Methyl palmitate, Dimethyl suberate	<b>Plant genomes</b>
<b>Source</b>	Pueraria lobata [TAX:3893]	
<b>Remark</b>	Therapeutic category: 5100	
<b>Comment</b>	Pueraria root	
<b>Other DBs</b>	PubChem: 47208344	
<b>LinkDB</b>	All DBs	

Formula: combination of crude drugs  
**Kakkonto**  
 DRUG: D06698

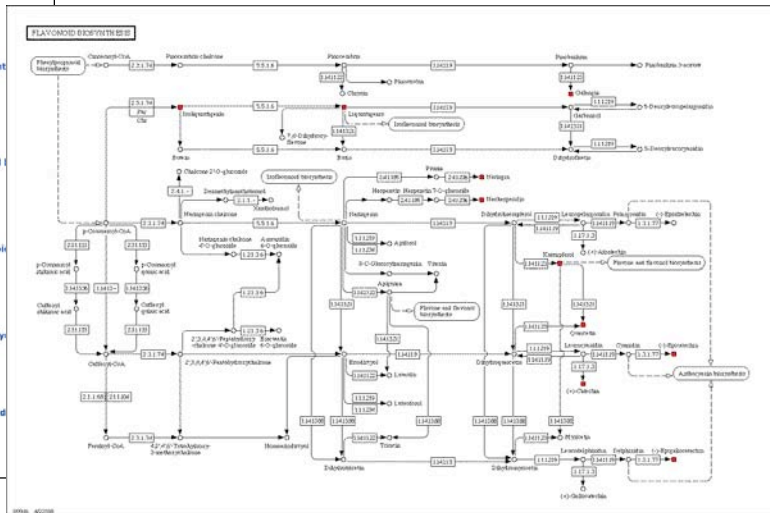
- ▼ Formulas
  - ▶ Formulas for relieving the exterior
    - ▶ **Formulas for relieving the exterior with pungent warmth**
      - D06698 Kakkonto extract (JP15); Kakkonto
      - D06928 Kakkontokassenkyushin'i
      - D06939 Keishikagoito
      - D06940 Keishikakakkonto
      - D06941 Keishikakobokukuyoninto
      - D06944 Keishikajutsubuto
      - D06947 Keishito
      - D06953 Keimakakuhanto
      - D06954 Kosogan
      - D06987 Shoseryuto
      - D07042 Maoto
      - D07045 Makuyokukanto
    - ▶ Formulas for relieving the exterior with pungent coolness
    - ▶ Formulas for relieving the exterior by supporting the normal
    - ▶ Purgative formulas
    - ▶ Formulas for mediation

Entry	D06698	Formula Drug
<b>Name</b>	Kakkonto extract (JP15); Kakkonto	
<b>Component</b>	Pueraria root [DR:D06693], Ephedra herb [DR:D06791], Jujube [DR:D06758], Peony root [DR:D06739 D06850], Glycyrrhiza [DR:D04365 D06828], Ginger [DR:D06744 D06852], Cinnamon bark [DR:D06712 D06841] or Cinnamon quassia [DR:D06898]	
<b>Activity</b>	Cold; Coryza	
<b>Remark</b>	Therapeutic category: 5200	
<b>Other DBs</b>	PubChem: 47208349	
<b>LinkDB</b>	All DBs	

# KEGG pathway mapping of chemical components in herbal medicine

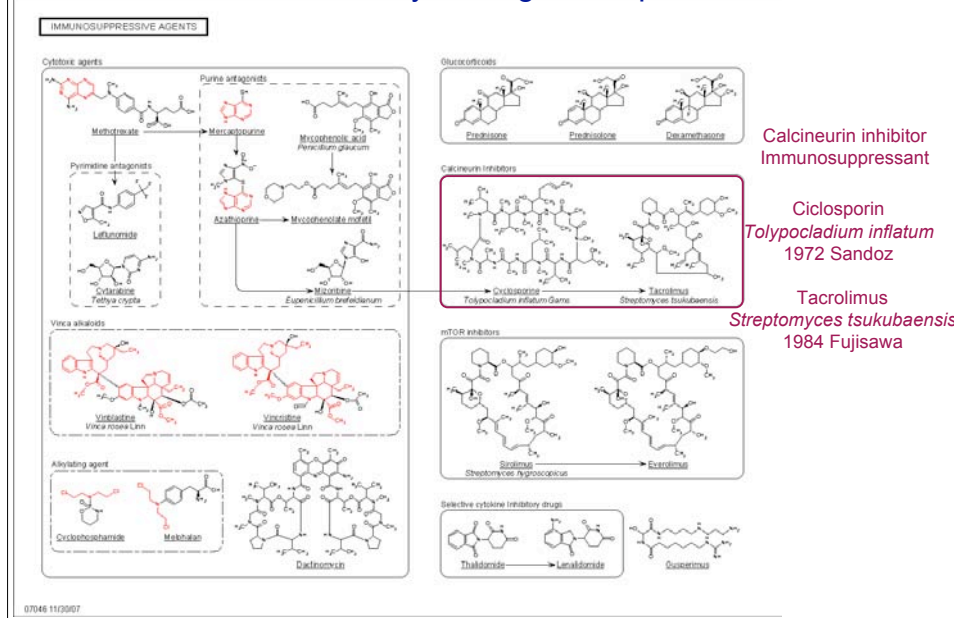
## Pathway Search Result

- **map00941 Flavonoid biosynthesis**
  - C00389
  - C03963
  - C05552
  - C06850
  - C07217
  - C09742
  - C09789
  - C09806
  - C10044
  - C12136
- **map00950 Alkaloid biosynthesis**
  - C00757
  - C01795
  - C02205
  - C02880
  - C04118
  - C05116
  - C05189
  - C05263
  - C05315
- **map00902 Monoterpenoid biosynthesis**
  - C00460
  - C00843
  - C04074
  - C06076
  - C09844
  - C09883
  - C11951
  - C11952
- **map00940 Phenylpropanoid biosynthesis**
  - C00911
  - C00952
  - C03003
  - C04984
  - C05267
  - C01533
  - C01752
  - C05851
- **map00061 Fatty acid biosynthesis**
  - C00249
  - C00712
  - C01510
  - C01571
  - C02639
  - C04623
  - C08362
- **map00281 Geraniol degradation**
  - C01499
  - C01900
  - C07847
  - C09848
  - C09851
  - C11386

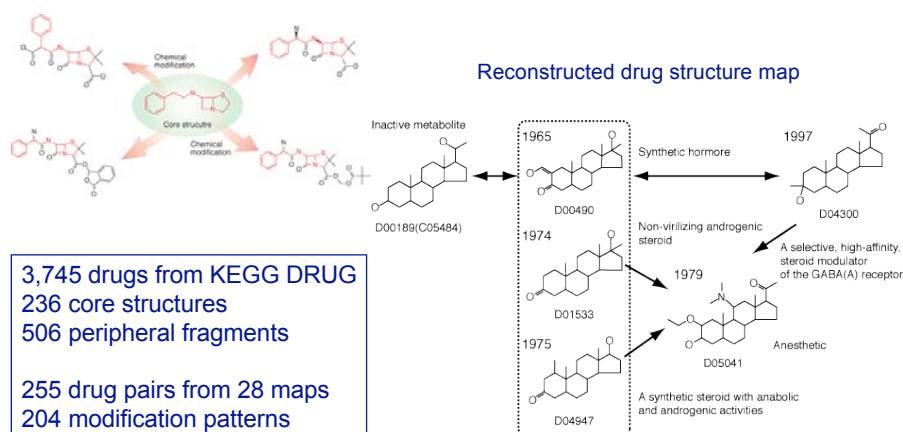




## KEGG DRUG structure maps for the history of drug development

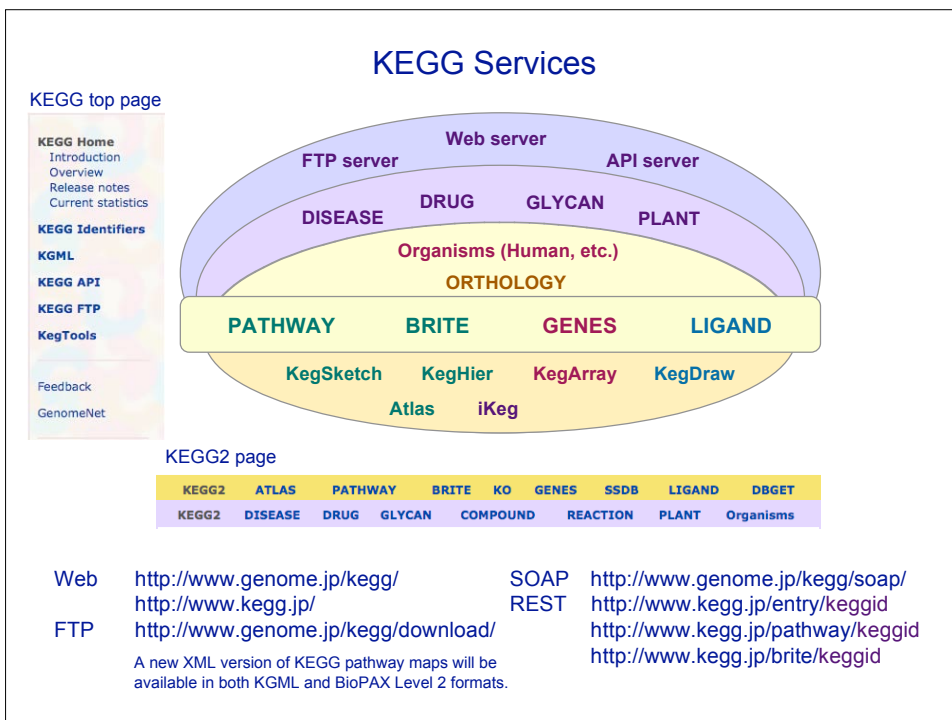


## Chemical modification patterns in drug development Towards understanding chemical architecture of marketed drugs



Shigemizu, D., Araki, M., Okuda, S., Goto, S., and Kanehisa, M.; Extraction and analysis of chemical modification patterns in drug development. *J. Chem. Info. Model.* 49, 1122-1129 (2009).

# Summary

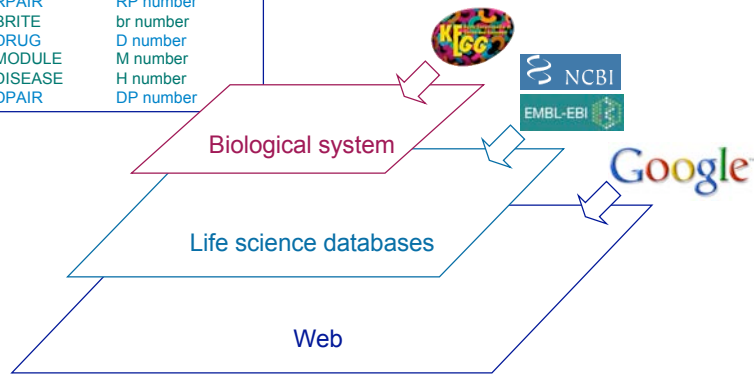


## KEGG Integrated with Other Resources

Namespace for KEGG objects

Release	Database	Object identifier
1995	KEGG PATHWAY	map number
	KEGG GENES	locus_tag / GeneID
	KEGG ENZYME	EC number
2000	KEGG COMPOUND	C number
	KEGG GENOME	organism code
2001	KEGG REACTION	R number
2002	KEGG ORTHOLOGY	K number
2003	KEGG GLYCAN	G number
2004	KEGG RPAIR	RP number
2005	KEGG BRITE	br number
	KEGG DRUG	D number
2007	KEGG MODULE	M number
	KEGG DISEASE	H number
2009	KEGG DPAIR	DP number

- KEGG is a computer representation of the biological systems.
- KEGG objects (database entries) are highly integrated.
- KEGG objects are linked to/from major life science databases.
- KEGG objects are part of the Web; they can be found by Web search engines.



## KEGG for biological interpretation of large-scale datasets

### Genes and proteins

Genome sequences  
Metagenomics data  
EST data  
Gene expression data  
etc.

### Identical K numbers

Pathway mapping  
BRITE mapping  
Structure mapping

KEGG PATHWAY  
Pathway maps  
Global maps  
Structure maps  
KEGG BRITE  
Functional hierarchies  
KEGG MODULE  
Gene/molecule lists  
KEGG DISEASE  
Gene/molecule lists

### Chemical substances

Metabolomics data  
Glycomics data  
etc.

### Identical C/D/G numbers

Structural similarity

Reference knowledge base  
of systemic biological functions

### Large-scale data sets

### Mapping as a set operation

### Binary relations

Protein-protein interactions  
Protein-small molecule interactions  
etc.

KEGG PATHWAY  
Pathway maps  
Global maps

Extended  
interaction  
network

## Changing roles of bioinformatics: molecules to molecular systems

---

### More divergent datasets

- Biopolymers: DNA, RNA, protein, glycan, lipid, polyketide, nonribosomal peptide, etc.
- Small molecules: metabolite, environmental compound, drug, etc.
- Molecular systems: pathway, complex, ontology, etc.

### More complex data types

- Biopolymers: linear sequence, circular sequence, tree structure, 3D structure
- Small molecules: chemical structure, 3D structure
- Molecular systems: graph, tree, list, etc.





### Hierarchy of data

- atom - molecule - molecular system
- monomer - biopolymer
- molecule - molecular interaction - molecular system
- cell - cellular interaction - cellular system

## Changing roles of bioinformatics: basic research to practical values

---

### Diseases viewed as perturbed states of molecular systems

- Capturing knowledge on molecular systems both in normal and perturbed states  KEGG PATHWAY  
KEGG DISEASE
- Capturing knowledge on drugs as perturbants to molecular systems  KEGG DRUG
- Capturing knowledge on natural product biosynthetic pathways and human-made chemical structure transformations  KEGG RPAIR  
KEGG DPAIR
- Generating knowledge from genome sequencing and other high-throughput data  KEGG ORTHOLOGY
- Knowledge based analysis of human diseases
- Drug discovery from the genomes of plants and microorganisms