

# Einführung zum Schwerpunktthema Datenbanken

Das Internet bietet Zugang zu einem einzigartigen Wissensfundus, der das gezielte Nachschlagen von Informationen in Lehrbüchern und gedruckten Fachzeitschriften zunehmend ersetzt. „Neue“ wissenschaftliche Erkenntnisse sind dank des rasanten Wissenszugewinns in der Genetik zunehmend bereits dann überholt, wenn die zugehörigen wissenschaftlichen Artikel im Druck erscheinen. Und „Online Supplementary Information“ bietet vielfach die einzige Möglichkeit, genetische Daten in ihrer Komplexität und Vernetzung mit anderen Datensätzen adäquat abzubilden. Darüber hinaus scheitern Printmedien bereits an der Darstellung der schiereren Menge der durch neue genetische Technologien erhobenen genetischen Daten: Als im Jahr 2000 die erste Draft-Sequenz des menschlichen Genoms verkündet wurde, waren Sequenzinformationen einer Größenordnung von acht Milliarden Basenpaaren (bp) in drei öffentlich zugänglichen und regelmäßig aufeinander abgestimmten Genom-Sequenzdatenbanken hinterlegt: GenBank am NCBI, dem US National Center for Biotechnology Information; RIKEN, der DNA-Datenbank Japans; sowie der Nucleotid-Sequenzdatenbank der Europäischen Molekularbiologischen Laboratorien (EMBL). Aktuell, nur 10 Jahre später, ist die Datenmenge der drei erwähnten Datenbanken auf 270 Milliarden bp angewachsen, und die Datenmenge verdoppelt sich etwa alle 18 Monate. Die Menge der gegenwärtig auf öffentlichen Datenservern hinterlegten Sequenz-Rohdaten wird in der Größenordnung von  $25 \times 10^{12}$  bp angege-

ben, von denen bereits jetzt der Großteil „new-generation“ Sequenzdaten individueller menschlicher Genome ausmachen [Nature Vol. 464, 01.04.2010]. Und die Speicherung von Datenmengen im Petabyte-Bereich ( $10^{15}$ ), wie sie bei der kompletten Sequenzierung ganzer Patientenkollektive erwartet werden, sind von öffentlich-geförderten Konsortien bereits nicht mehr tragbar, so dass das Datenmanagement kommerziellen Anbietern übertragen wird.

Ist die Humangenetik also einerseits durch riesige, bislang weitgehend nur unzulänglich interpretierbare Sequenz-Datenmengen herausgefordert, hat sie andererseits mit der erheblichen klinischen und biologischen Vielfalt genetischer Krankheitsbilder zu kämpfen, deren komplexe Phänotypen sich einer standardisierten Beschreibung häufig entziehen. Die Anzahl individueller Datenbanken und Algorithmen, die sich beiden Ebenen und ihren Zusammenhängen widmen, ist kaum mehr überschaubar. Durch den Alltag am Computer hat sich die Arbeitsweise in der humangenetischen Beratung und Diagnostik grundlegend verändert: Die Anwendung bioinformatischer Datenbanken reicht von primär krankheitsbezogenen Literaturrecherchen über die Interpretation eigener Laborbefunde bis hin zu komplexen statistischen Analysen in Populationsdatensätzen. Der beschränkte Rahmen einer Ausgabe der Zeitschrift „medizinische Genetik“ zum Schwerpunktthema Datenbanken kann daher keinen systematischen und erst recht keinen umfassenden Zugang

zu diesem Themengebiet gewährleisten. Stattdessen haben sich die wissenschaftlichen Koordinatoren entschlossen, aktuelle Trends anhand einiger ausgewählter Beiträge von Arbeitsgruppen in Deutschland darzustellen. Darüber hinaus soll eine Toolbox mit einer reichhaltigen Linksammlung das Stöbern im Internet auf eigene Initiative stimulieren.

Bei einer anonymisierten Umfrage im Kollegenkreis des Instituts der Koordinatoren gab jeder der Teilnehmenden zu, humangenetisch-relevante Informationen mit großer Regelmäßigkeit zu „googlen“. Auch wenn mehrere Studien zeigen, dass gezielte Google-Anfragen gerade bei seltenen Erkrankungen durchaus zur richtigen Diagnosestellung führen können [z.B. Tang und Ng, 2006 BMJ, 333:1143], basiert das Suchergebnis üblicherweise auf der Strategie, wo und wie man sucht – noch fragen Computer im Allgemeinen nicht nach. Während manchmal ein einzelnes „richtiges“ Schlagwort zielsicher zum gewünschten Ergebnis führt, verliert sich der typische User regelmäßig in den Weiten des Netzes.

Im ersten Artikel dieses Hefts stellen Rommel und Kollegen daher eine Option zum „Wo“ vor. **Orphanet** ist eine europäische Datenbank für seltene Krankheiten, zu der Partner aus über 30 Ländern Krankheits- und Patienteninformationen beitragen. Ärzte, Forscher und vor allem Patienten selbst können sich in verschiedenen Sprachen nicht nur über Details einer Vielzahl seltener Krankheitsbilder informieren, sondern werden darüber hinaus auf diagnostische Tests, Spe-

zialambulanzen, Selbsthilfegruppen oder verfügbare Therapieoptionen hingewiesen. Hinsichtlich des „Wie“ dürfte sich die Google-Trefferquote bei Verwendung eines standardisierten Vokabulars zur präzisen Beschreibung humangenetischer Sachverhalte wesentlich erhöhen. Die systematische Erfassung eines komplexen Krankheits-Phänotyps stellt generell jedoch eine ungleich größere Herausforderung dar, als die Erstellung primärer Sequenzinformation. Über einen ersten Schritt dazu – der Erstellung einer gemeinsamen Sprache, der **Human Phenotype Ontology** – berichten Doelken und Kollegen in ihrem Artikel.

**Biomaterial-Datenbanken** – strukturierte, qualitativ hochwertige Sammlungen nicht nur virtueller Datensätze, sondern auch biologischer Materialien, möglichst gemeinsam mit Detailinformationen zum Phänotyp der zugehörigen Individuen, sind eine wichtige Quelle für die populationsgenetische und genetisch-epidemiologische Forschung. Krawczak und Kollegen geben in ihrem Beitrag einen Überblick und diskutieren die rechtlichen, ethischen und strukturellen Herausforderungen von Biobanken. Die langfristige Verfügbarkeit und die Bearbeitung auch zukünftiger Forschungsfragen erfordern spezifische datenschutzrechtliche Erwägungen, die im Voraus sicherstellen müssen, dass die Interessen aller Beteiligten gewahrt bleiben. Aktuell sind Biobanken die Grundlage vieler genom-weiter Assoziationsstudien (GWAS), deren unmittelbare Bedeutung für die genetische Beratung und Diagnostik gegenwärtig lebhaft diskutiert wird. So beziehen sich kommerzielle Anbieter bereits jetzt auf GWAS-Resultate, um Krankheitsprädispositionen zu diagnostizieren und Strategien zur Risikoreduktion zu empfehlen. Es ist daher davon auszugehen, dass die Interpretation von GWAS-Ergebnissen für die genetische Beratung zunehmende Relevanz erhalten wird. Eine systematische Auswertung von GWAS-Studien durch **GWAS-Metaanalysen** erlaubt nicht nur eine Objektivierung der tatsächlichen Signifikanz einzelner Loci, sondern könnte zur Verschmälerung der Lücke zwischen Grundlagenwissenschaft und klinischer Anwendung beitragen. Dies stellen Lill und Bertam in ihrem Artikel etwa am Beispiel der

Alzheimerkrankheit auf beeindruckende Weise dar.

Es ist davon auszugehen, dass die Diskrepanz zwischen erhobener Datenmenge und ihrer Interpretierbarkeit durch die nächste Stufe der technologischen Entwicklung, der schnellen und kostengünstigen Hochdurchsatz-Sequenzierung individueller Genome, sogar noch zunehmen wird. Stütz und Korbel beschäftigen sich in ihrem Artikel zum **1000-Genome-Project** schwerpunktmäßig mit der informationstechnologischen Infrastruktur zur Archivierung riesiger Datenmengen und deren Bearbeitung, die bei der vollständigen Genomsequenzierung von bis zu 1000 Individuen anfallen – ein Projekt, das bis 2011 weitgehend abgeschlossen sein soll. Die Autoren prognostizieren, dass sich die Arbeit der Humangenetiker durch zukünftig zu erwartende sehr viel umfangreichere Datenanalysen noch mehr als bisher wandeln wird.

Mögliche Szenarien, mit denen sich der Genetische Berater der (nahen) Zukunft auseinandersetzen zu haben dürfte, sind dem Beitrag von Krawitz zu entnehmen, in dessen Mittelpunkt die Bedeutung von „**Personal Genomics**“ für die Prognose und zukünftige Behandlung des individualisierten Patienten steht. Während eine maßgeschneiderte Therapie von Volkskrankheiten gegenwärtig wohl noch als Zukunftsmusik erachtet werden kann, zeigen einige Beispiele aus der Pharmakogenomik, dass bestimmte Genvarianten große Effekte auf Medikamentenwirkungen haben können, so dass individualisierte Genominformation durchaus von Nutzen sein kann.

Während das typische humangenetische Institut der Gegenwart von solchen Fragestellungen zumindest wohl noch eine Zeit lang verschont bleiben wird, gehört das informationstechnologische Management erhobener und gespeicherter Daten in jeder genetischen Beratungsstelle und Labor zum Alltag. Die abschließenden Artikel von Schlott/Schröck und Schröder/Müller-Reible zeigen am Beispiel zweier ausgewählter **Patienten- und Labordaten-Managementsysteme** wie Patienten-, Familien- und Labordaten in humangenetischen Praxen und Einrichtungen verwaltet werden können. Beide vorgestellten Systeme wurden im aka-

demischen Umfeld aus der Praxis heraus zum Management der eigenen Daten entwickelt und sind daher direkt auf die Erfordernisse der gemeinsamen Verwaltung von Prozessen wie Genetische Beratung, Probeneingang, Analyseverlauf, Ergebnisdokumentation, Befunderstellung und Abrechnung zugeschnitten. Bei konsequenter Umsetzung der verfügbaren Möglichkeiten lässt sich durch den Einsatz derartiger Programme sehr viel doppelte Datenführung vermeiden, was zu einem weitgehend papierlosen und insgesamt sehr viel effektiveren Arbeiten beitragen sollte.

### Korrespondenzadresse

**H. Runz, Ch. Fischer**  
Institut für Humangenetik  
INF 366  
Universität Heidelberg  
69120 Heidelberg  
heiko.runz@med.uni-heidelberg.de

# Datenbanken-Toolbox

(Stand: 04/2010)

Die nachfolgende Toolbox listet Links zu einer Auswahl an online frei verfügbaren (open-source) Datenbanken, die Informationen und Analyse-Tools für häufige Fragestellungen in der Humangenetik anbieten. Die Auswahl erhebt keinen Anspruch auf Vollständigkeit, und im schnelllebigen World Wide Web kann eine Garantie weder für die Qualität des Seiteninhalts, noch für die Aktualität der aufgelisteten Internet-Adressen übernommen werden. Die Liste richtet sich an den „typischen“

humangenetischen User, der als Anwender Informationen aus bestehenden Datensätzen zusammentragen oder eigene kleinere Datensätze (z.B. mögliche funktionelle Konsequenzen einer bestimmten DNA-Sequenzvariante in einem einzelnen Gen) selbst analysieren möchte. Von der Auflistung der zahlreichen und häufig exzellenten Genlokus-spezifischen Datenbanken wurde ebenso abgesehen wie von der Anführung von Tools (z.B. zum Sequenz-Alignment) die wesentliche bioin-

formatische Vorkenntnisse voraussetzen. Für die Mehrzahl der Seiten sind ausführliche Online-Tutorials verfügbar, die eine schnelle Einführung in deren adäquate Nutzung erlauben. Die aufgeführte Liste ist unter der URL <http://www.gfhev.de/de/links/fachinformationen.htm> auch online verfügbar. Verbesserungsvorschläge und Tipps zur Aktualisierung und Ergänzung weiterer relevanter Links sind per Email an [organisation@gfhev.de](mailto:organisation@gfhev.de) herzlich willkommen. Viel Erfolg beim Surfen!

## 1. Literatur- und Meta-Suchmaschinen:

|                |   |
|----------------|---|
| Google         | <a href="http://www.google.com">http://www.google.com</a>                           |
| Google Scholar | <a href="http://scholar.google.com">http://scholar.google.com</a>                   |
| Google Books   | <a href="http://books.google.com">http://books.google.com</a>                       |
| NCBI PubMed    | <a href="http://www.ncbi.nlm.nih.gov/PubMed">http://www.ncbi.nlm.nih.gov/PubMed</a> |
| Gopubmed       | <a href="http://www.gopubmed.com">http://www.gopubmed.com</a>                       |

## 2. Datenbanken mit Krankheitsinformationen/ Relevanz für die klinische Genetik:

|                       |   |
|-----------------------|---|
| OMIM                  | <a href="http://www.ncbi.nlm.nih.gov/omim">http://www.ncbi.nlm.nih.gov/omim</a>                                   |
| GeneReviews           | <a href="http://www.ncbi.nlm.nih.gov/sites/GeneTests">http://www.ncbi.nlm.nih.gov/sites/GeneTests</a>             |
| Orphanet              | <a href="http://www.orpha.net/consor/cgi-bin/index.php">http://www.orpha.net/consor/cgi-bin/index.php</a>         |
| NORD                  | <a href="http://www.rarediseases.org">http://www.rarediseases.org</a>   |
| EURORDIS              | <a href="http://www.eurordis.org">http://www.eurordis.org</a>   |
| Arzneimittelwirkungen | <a href="http://www.arzneimittel-in-der-Schwangerschaft.de">http://www.arzneimittel-in-der-Schwangerschaft.de</a> |

## 3. Datenbanken mit Relevanz für die klassische Zytogenetik / FISH:

|                                |   |
|--------------------------------|---|
| Mendelian Cytogenetics Network | <a href="http://www.mcndb.org">http://www.mcndb.org</a>   |
| ENSEMBL Cytoview               | <a href="http://www.ensembl.org/Homo_sapiens/cytoview">http://www.ensembl.org/Homo_sapiens/cytoview</a>             |
| Zytogenetik-Atlas (Onkologie)  | <a href="http://atlasgeneticsoncol.org/">http://atlasgeneticsoncol.org/</a>   |
| Mitelman-db (Onkologie)        | <a href="http://cgap.nci.nih.gov/Chromosomes/Mitelman">http://cgap.nci.nih.gov/Chromosomes/Mitelman</a>             |
| SKY/M-FISH, CGH-db (Onkologie) | <a href="http://www.ncbi.nlm.nih.gov/projects/sky">http://www.ncbi.nlm.nih.gov/projects/sky</a>                     |
| Progenetix (Onkologie)         | <a href="http://www.progenetix.net/progenetix/index.html">http://www.progenetix.net/progenetix/index.html</a>       |
| Array-CGH Tumor-db (Onkologie) | <a href="http://amba.charite.de/~ksch/cghdatabase/index.htm">http://amba.charite.de/~ksch/cghdatabase/index.htm</a> |

## 4. Humanes Genom Annotations-Browser / Analyse-Tools:

|                        |   |
|------------------------|---|
| Ensembl                | <a href="http://www.ensembl.org/index.html">http://www.ensembl.org/index.html</a>                                       |
| UCSC Genome Browser    | <a href="http://genome.ucsc.edu/">http://genome.ucsc.edu/</a>   |
| NCBI Map Viewer        | <a href="http://www.ncbi.nlm.nih.gov/projects/mapview">http://www.ncbi.nlm.nih.gov/projects/mapview</a>                 |
| Galaxy Genome Software | <a href="http://bitbucket.org/galaxy/galaxy-central/wiki/Home">http://bitbucket.org/galaxy/galaxy-central/wiki/Home</a> |
| Biomart                | <a href="http://www.biomart.org">http://www.biomart.org</a>   |
| DAS                    | <a href="http://biodas.org">http://biodas.org</a>   |

## 5. Gen-basierte Datenbanken/ Meta-Suchmaschinen:

|                         |   |
|-------------------------|---|
| HGNC                    | <a href="http://www.genenames.org">http://www.genenames.org</a>   |
| NCBI Entrez Gene        | <a href="http://www.ncbi.nlm.nih.gov/entrez/query.fcgi?db=gene">http://www.ncbi.nlm.nih.gov/entrez/query.fcgi?db=gene</a> |
| GeneCards               | <a href="http://www.genecards.org">http://www.genecards.org</a>   |
| Bioinformatic Harvester | <a href="http://harvester.fzk.de">http://harvester.fzk.de</a>   |
| Genedistiller           | <a href="http://www.genedistiller.org">http://www.genedistiller.org</a>   |
| iHOP                    | <a href="http://www.ihop-net.org/UniPub/iHOP/">http://www.ihop-net.org/UniPub/iHOP/</a>                                   |

**6. Genom-Variations Datenbanken / Genotyp-Phänotyp-Korrelation:**

|                                       |   |
|---------------------------------------|---|
| NCBI dbSNP                            | <a href="http://www.ncbi.nlm.nih.gov/SNP">http://www.ncbi.nlm.nih.gov/SNP</a>                                 |
| Hapmap                                | <a href="http://hapmap.ncbi.nlm.nih.gov/">http://hapmap.ncbi.nlm.nih.gov/</a>                                 |
| Database of Genomic Variants          | <a href="http://projects.tcag.ca/variation">http://projects.tcag.ca/variation</a>                             |
| Decipher                              | <a href="https://decipher.sanger.ac.uk/application">https://decipher.sanger.ac.uk/application</a>             |
| Segmental duplication db              | <a href="http://humanparalogy.gs.washington.edu">http://humanparalogy.gs.washington.edu</a>                   |
| NCBI dbGAP                            | <a href="http://www.ncbi.nlm.nih.gov/sites/entrez?db=gap">http://www.ncbi.nlm.nih.gov/sites/entrez?db=gap</a> |
| NCBI Genetic Association db           | <a href="http://geneticassociationdb.nih.gov">http://geneticassociationdb.nih.gov</a>                         |
| Gen2Phen Knowledge Center             | <a href="http://www.gen2phen.org/">http://www.gen2phen.org/</a>   |
| European Genome-Phenome Archive       | <a href="http://www.ebi.ac.uk/ega">http://www.ebi.ac.uk/ega</a>   |
| Welcome Trust Case Control Consortium | <a href="http://www.wtcc.org.uk">http://www.wtcc.org.uk</a>   |
| Human Epigenome Project               | <a href="http://www.epigenome.org">http://www.epigenome.org</a>   |
| Human Gene Mutation Database          | <a href="http://www.hgmd.cf.ac.uk/ac/index.php">http://www.hgmd.cf.ac.uk/ac/index.php</a>                     |
| HGVS (Nomenklatur)                    | <a href="http://www.hgvs.org/mutnomen">http://www.hgvs.org/mutnomen</a>                                       |
| Mutalyzer (Nomenklatur)               | <a href="http://www.LOVD.nl/mutalyzer">http://www.LOVD.nl/mutalyzer</a>                                       |
| Locus reference Genomic (Nomenklatur) | <a href="http://www.lrg-sequence.org">http://www.lrg-sequence.org</a>   |

**7. Gensequenzanalyse-Tools:**

**7.1. Sequenzvarianten-Interpretations-Tools:**

|                 |   |
|-----------------|---|
| Polyphen2       | <a href="http://genetics.bwh.harvard.edu/pph2">http://genetics.bwh.harvard.edu/pph2</a>                                 |
| SIFT            | <a href="http://sift.jcvi.org/">http://sift.jcvi.org/</a>   |
| Mutation Taster | <a href="http://neurocore.charite.de/MutationTaster/">http://neurocore.charite.de/MutationTaster/</a>                   |
| SNPS3D          | <a href="http://www.snps3d.org">http://www.snps3d.org</a>   |
| SNAP            | <a href="http://www.rostlab.org/services/SNAP">http://www.rostlab.org/services/SNAP</a>                                 |
| Pmut            | <a href="http://mmb2.pcb.ub.es:8080/PMut">http://mmb2.pcb.ub.es:8080/PMut</a>   |
| AGVGD           | <a href="http://agvgd.iarc.fr/index.php">http://agvgd.iarc.fr/index.php</a>   |
| Modbase/LS-SNP  | <a href="http://modbase.compbio.ucsf.edu/LS-SNP//About.html">http://modbase.compbio.ucsf.edu/LS-SNP//About.html</a>     |
| SNPeffect       | <a href="http://snpeffect.vib.be/index.php">http://snpeffect.vib.be/index.php</a>                                       |
| TopoSNP         | <a href="http://gila-fw.bioengr.uic.edu/snp/toposnp">http://gila-fw.bioengr.uic.edu/snp/toposnp</a>                     |
| MutDB           | <a href="http://www.mutdb.org">http://www.mutdb.org</a>   |
| pupaSNP         | <a href="http://pupasnp.org">http://pupasnp.org</a>   |
| FastSNP         | <a href="http://fastsnp.ibms.sinica.edu.tw">http://fastsnp.ibms.sinica.edu.tw</a>                                       |
| ssahaSNP        | <a href="http://www.sanger.ac.uk/resources/software/ssahasnp/">http://www.sanger.ac.uk/resources/software/ssahasnp/</a> |

**7.2. SpliceSite-Vorhersage-Tools:**

|  |   |
|--|---|
| HSF                                      | <a href="http://www.umd.be/HSF/">http://www.umd.be/HSF/</a>   |
| ACESCAN                                  | <a href="http://genes.mit.edu/acescan/">http://genes.mit.edu/acescan/</a>                                 |
| Splice Site Prediction (D. melanogaster) | <a href="http://www.fruitfly.org/seq_tools/splice.html">http://www.fruitfly.org/seq_tools/splice.html</a> |

**7.3. DNA-Sequenzmotiv-Vorhersage-Tools:**

|          |   |
|----------|---|
| TFSEARCH | <a href="http://www.cbrc.jp/research/db/TFSEARCH.html">http://www.cbrc.jp/research/db/TFSEARCH.html</a> |
| JASPAR   | <a href="http://jaspar.cgb.ki.se">http://jaspar.cgb.ki.se</a>   |
| MEME     | <a href="http://meme.sdsc.edu/meme">http://meme.sdsc.edu/meme</a>                                       |
| TRAWLER  | <a href="http://ani.embl.de/trawler">http://ani.embl.de/trawler</a>                                     |
| WEEDER   | <a href="http://159.149.109.9/modtools/">http://159.149.109.9/modtools/</a>                             |

**8. ausgewählte Gen-/ Proteinfunktions Datenbanken:**

|                       |   |
|-----------------------|---|
| GeneOntology          | <a href="http://www.geneontology.org">http://www.geneontology.org</a>                       |
| SRS 3d                | <a href="http://srs3d.org">http://srs3d.org</a>   |
| Pfam                  | <a href="http://pfam.sanger.ac.uk/">http://pfam.sanger.ac.uk/</a>                           |
| Expasy                | <a href="http://www.expasy.ch">http://www.expasy.ch</a>                                     |
| PDB                   | <a href="http://www.pdb.org">http://www.pdb.org</a>   |
| FOLDX                 | <a href="http://foldx.crg.es/about.jsp">http://foldx.crg.es/about.jsp</a>                   |
| TOPO2                 | <a href="http://www.sacs.ucsf.edu/TOPO2">http://www.sacs.ucsf.edu/TOPO2</a>                 |
| KEGG pathway analysis | <a href="http://www.genome.ad.jp/kegg">http://www.genome.ad.jp/kegg</a>                     |
| Proteinatlas          | <a href="http://www.proteinatlas.org">http://www.proteinatlas.org</a>                       |
| 4Dxpress              | <a href="http://4dx.embl.de/4DXpress/welcome.do">http://4dx.embl.de/4DXpress/welcome.do</a> |
| Mitochcek             | <a href="http://www.mitochcek.org/">http://www.mitochcek.org/</a>                           |

**Sammlung von Programmen zur statistischen Analyse von genetischen Daten**

incl. Stammbaumzeichenprogramme: <http://linkage.rockefeller.edu/soft/>