

Unexplored opportunities in the druggable human genome

The IDG Knowledge Management Center

Much of biomedical research and the development of therapeutics is focused on a small fraction of the human genome, ignoring many disease-relevant proteins and the associated scientific and commercial opportunities. The National Institutes of Health (NIH) Illuminating the Druggable Genome programme aims to catalyse research around understudied targets. As part of this effort, we categorized proteins based on available data, information and knowledge from various sources. Furthermore, we evaluated the commercial focus on therapeutic targets and

21%

associated disease areas, by processing global drug sales data for a 5-year period. Taken together, our quantification of available data and knowledge, as well as the financial analyses based on global drug sales and NIH funding for that same period, highlight a bias towards well-described proteins, particularly drug targets. Our analysis identifies a knowledge deficit, not just concerning experimental data, but also with respect to patents and publications. We believe this points towards novel research and commercial opportunities for drug discovery and development.



Target development levels

with the T_{clin} category expanded.

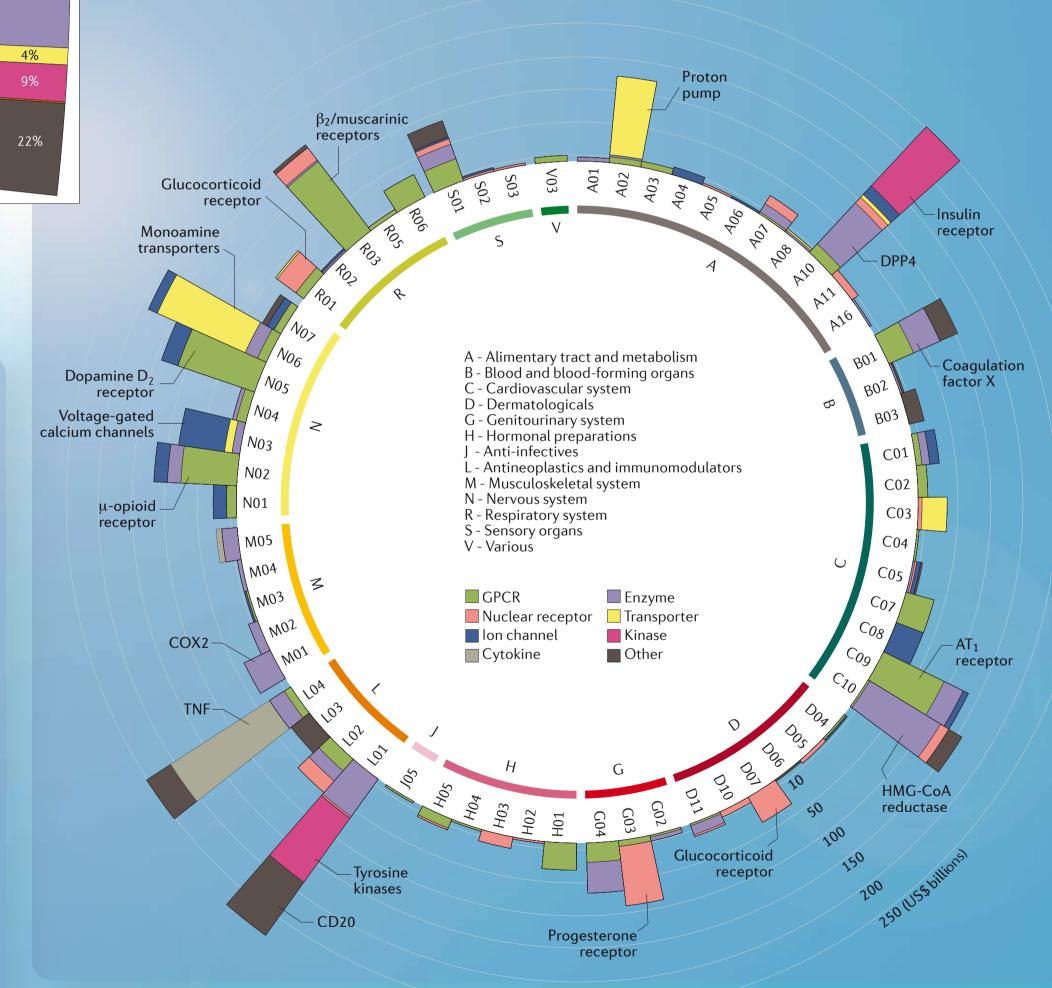
Most current protein classification schemes are based on structural and functional criteria. For target prioritization and therapeutic development it is useful to understand how much and what types of data are available for a given protein, thereby highlighting well-studied and understudied targets. Thus a classification scheme based on the amount and types of knowledge that are available can lay the foundation for rational target prioritization. To address this, we developed the target development level (TDL) classification scheme (see footnote*), which uses three criteria to classify a protein into one of four categories, according to the following ordering. Proteins annotated as drug targets by mode of action¹ were categorized as T_{atta}. Proteins for which small-molecule modulators (with reported bioactivity above specific cut-off values; see footnote[‡]) are known were categorized as T_{stor}. Proteins categorized as T_{stor} either have confirmed OMIM disease phenotypes or Gene Ontology (GO) annotations² based on experimental evidence, or meet two of the following three conditions: a fractional count of PubMed abstracts mentioning the protein³ above five, three or more Gene Reference Into Function (RIF) annotations⁴, or more than fifty antibodies according to Antibodypedia⁵. The remaining proteins 20% are classified as T_{dark} and are the proteins for which there is the least knowledge and the lowest number of molecular probes available. Moreover, research on these proteins receives relatively little funding, as indicated by the distribution of NIH R01 grants below. The figure to the immediate right shows each TDL category as percentage of the whole proteome in the inner ring, and percentages of each TDL category for selected major target families in the outer ring,

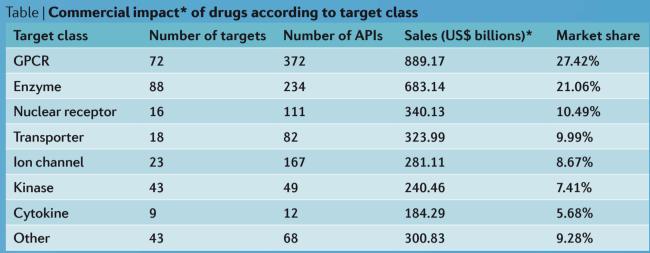
Our analysis of commercial activity on the right is based on global drug sales, which stem from only 3% of the proteome (T_{alia}). Our ability to specifically and selectively perturb proteins using small molecules is apparently limited to less than 10% of the human proteome. However, as many of these proteins interact and re-arrange in multiple ways, it is likely that the proportion of targets that can be perturbed by small molecules exceeds 10%.

A target-centric analysis of global drug sales data

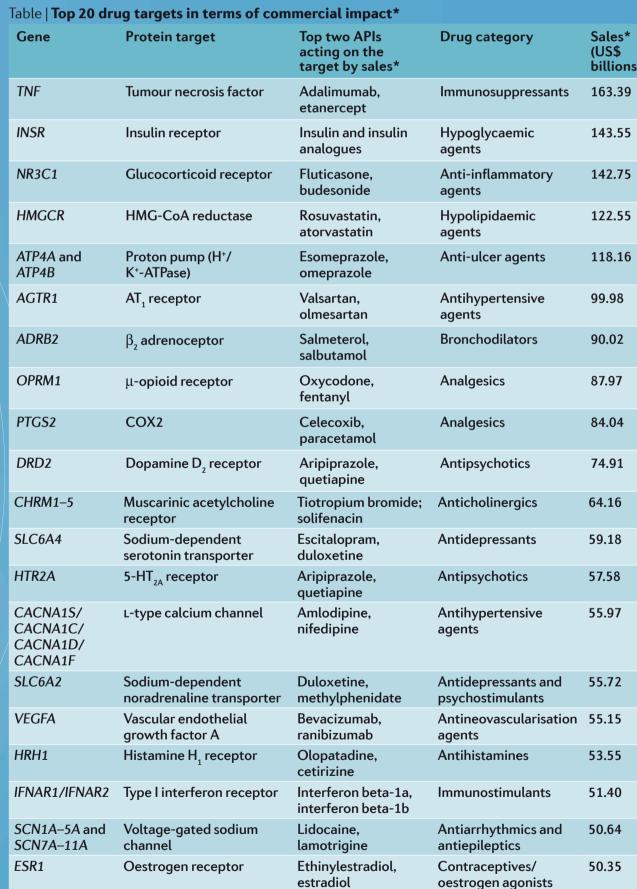
While evaluating the current level of knowledge about current and potential drug targets, we asked the question: "what are the most lucrative targets?" from a therapeutic perspective. To investigate this question, we analysed IMS Health data on drug sales from 75 countries, including Japan and those in Europe and North America, aggregated over a 5-year period (2011–2015). After excluding botanicals, traditional Chinese and homeopathic medicines and drugs perturbing non-human targets (such as some antiviral drugs), we identified 51,095 unique products. These were mapped to 1,069 active pharmaceutical ingredients (APIs) from DrugCentral¹⁰, corrected by number of APIs per product and by number of Talia targets per API. Sales organized by ATC level 2 codes¹¹ and by target class were normalized to percentage values, and depicted in a circular histogram (below); the value of the sales is indicated in the concentric circles around the histogram, with a scale in the bottom right. The top two ATC level 1 categories in terms of global sales of the associated APIs are 'antineoplastics and immunomodulators' followed by the 'nervous system'.

The commercial impact of the seven target classes discussed in this poster is summarized in the table on the top right. The most lucrative targets are G protein-coupled receptors (GPCRs), with ~27% of the overall sales. Enzymes (excluding kinases, which are considered separately given their importance), are the next class, followed by nuclear receptors and transporters. Sales of drugs that modulate these seven target classes, which include 269 T., targets and 1,027 APIs, reached nearly US\$3 trillion for the 2011–2015 period. Cytokines, the latest addition (in the early 2000s) to the list, are currently the most lucrative target. The top 20 T_{dis} protein targets, ranked by aggregated API sales data, are listed in the table on the bottom right.





*Sales are aggregated over the period 2011–2015.



*Sales are aggregated over the period 2011–2015. 5-HT, 5 hydroxytryptamine. AT_1 , angiotensin II type 1; COX2, cyclooxygenase 2; DPP4, dipeptidyl peptidase 4; HMG-CoA, 3-hydroxy-3-methylglutaryl-coenzyme A.

External validation of target development levels and the knowledge deficit

Different resources and data sets were used to examine the lack of annotation of proteins catalogued as T., compared with other proteins (see footnote*). From left to right, the first three normalized score distributions (on a green background) illustrate the criteria used to assign target development levels, namely PubMed abstracts, the Gene RIF counts and the antibody counts per protein.

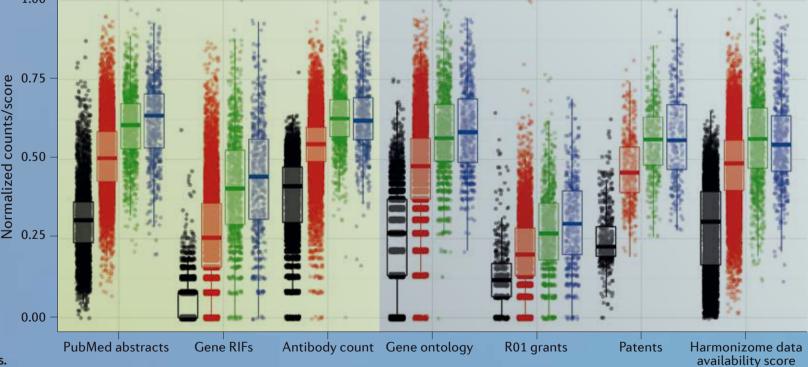
These criteria were externally validated on four other sets of data (on a blue background) as follows:

• 'Gene ontology terms' illustrates the distribution of Gene Ontology² annotation counts per protein, using data from UniProt⁶. More than 90% of T_{det} proteins, but less than 10% of T_{elia} proteins, have 10 Gene Ontology annotations or less.

• 'R01 grants' illustrates the distribution of text-mined NIH R01 grant counts detected for each protein, using data from NIH RePORTER'. Most T_{dark} proteins do not have research on them funded via the NIH R01 mechanism.

• 'Patents' illustrates the distribution of text-mined granted patents for each protein, using data from SureChEMBL⁸. T_{dest} proteins are included in a significantly lower number of patents compared with proteins in the other categories. • 'Harmonizome data availability score' illustrates experimental information density per protein, processed from 70 genomic data sets9. T. ... proteins appear to have

Through the IDG KMC project, various sources of information related to biological assays, human and mouse phenotypes, disease associations, expression data from transcriptomic and proteomic experiments, pathways and other processes and functions have been integrated and harmonized. There is an uneven distribution of information for targets in multiple types of information, including research literature, patent literature and experimental data. In particular, $\overset{\circ}{2}$ 0.25 there appears to be a substantial knowledge deficit for T_{dark} proteins, suggesting a considerable opportunity for novel discoveries, and ultimately new commercial opportunities, if further resources are directed towards investigating such proteins.



Transporter

Epigenetic

Transcription factor

Kinase

Other

Nuclear recepto

Ion channel

Enzyme

*The Target Central Resource Database, version 3.0.7 (September 2016), contains 20,186 human proteins (manually reviewed Uniprot entries only), which were classified into four target development levels (TDLs): T_{dot} (7,583 proteins), T_{bio} (10,759 proteins), T_{chem} (1,243 proteins), and T_{clin} (601 proteins). TDL levels were established by evaluating knowledge, as shown above. Each data point represents a protein, colour-coded by TDL; boxplots represent the interquartile range (25% to 75%) and the median for each knowledge availability metric. Raw counts and scores were first converted to log scale, then normalized to unit variance. kinases: ≤30 nM; G protein-coupled receptors: ≤100 nM; nuclear receptors: ≤100 nM; ion channels: $\leq 10 \,\mu\text{M}$; non-IDG family targets: $\leq 1 \,\mu\text{M}$. T_{chan} proteins have

bioactivities in ChEMBL¹² and DrugCentral and human curation for some targets.

1. Santos, R. et al. A comprehensive map of molecular drug targets. Nat. Rev. Drug Disc. doi: 10.1038/nrd.2016.230 (2016). 2. Ashburner, M. Gene ontology: tool for the unification of biology. The Gene Ontology Consortium. Nat. Genet. 25,

- 6. UniProt Knowledgebase: http://www.uniprot.org
- 7. NIH RePORTER portal: https://projectreporter.nih.gov/reporter.cfm
- 8. SureChEMBL portal: https://www.surechembl.org/search/
- 10. Ursu, O. et al. DrugCentral: online drug compendium. Nucl. Acids Res. doi: 10.1093/nar/gkw993; http://drugcentral.org 11. ATC classification system: http://www.whocc.no/atc_ddd_index

- Cristian Bologa, Jayme Holmes, Subramani Mani, Steve L. Mathias, Tudor I. Oprea, Larry A. Sklar, Oleg Ursu, Anna Waller, Jeremy J. Yang and Gergely Zahoranszky-Kohalmi* are at the University of New Mexico, MSC09 5025, One University of New Mexico, 700 Camino de Salud NE, Albuquerque,
- Lars Juhl Jensen and Søren Brunak are at the Novo Nordisk Foundation Center for Protein Research, University of Copenhagen, Blegdamsvej 3b,
- Avi Ma'ayan and Joel Dudley are at the Icahn School of Medicine at Mount Sinai, One Gustave L. Levy Place, Box 1603, New York, New York 10029, USA. formulations, representing aggregated 5-year sales for over 1,200 active pharmaceutical • John Overington[§], Anne Hersey, Anna Gaulton, Anneli Karlson[§] and George Papadatos[‡] are at the European Molecular Biology Laboratory, European
- Rajarshi Guha, Ajit Jadhav, Dac-Trung Nguyen, Anton Simeonov and Noel Southall are at NIH-NCATS, 9800 Medical Center Drive, Rockville,
- Stephan Schürer and Dušica Vidović are at Department of Pharmacology, Miller School of Medicine and Center for Computational Science,
- Allen Campbell and Christian Reich are at IMS Health, 2200 Bridge Parkway, Redwood City, California 94065, USA.

Bioinformatics Institute, Wellcome Genome Campus, Hinxton, Cambridge, CB10 1SD, UK.

Present addresses: *NIH-NCATS, 9800 Medical Center Drive, Rockville, Maryland 20852, USA. *GlaxoSmithKline, Gunnels Wood Road, Stevenage, Hertfordshire, SG1 2NY, UK. \$BenevolentAl, 40 Churchway, London, NW1 1LW, UK. Correspondence to T.I.O. toprea@salud.unm.edu

The authors are thankful to IMS Health for providing data for over 60,000 pharmaceutical ingredients from 75 countries. Funding for this work was provided by NIH grant U54CA189205 (to T.I.O., A.S. and L.A.S.) and by NIH grant U54CA189201 (to A.M. & J.D). The production and distribution of this poster was funded by the Illuminating the Druggable Genome Knowledge Management Center, NIH grant U54CA189205. The poster content is peer reviewed and the responsibility of Macmillan Publishers Limited.

Edited by Peter Kirkpatrick; copy-edited by Carrie Hardy. © 2016 Macmillan Publishers Limited. All rights reserved. http://www.nature.com/nrd/posters/druggablegenome/index.html

Illuminating the Druggable Genome Knowledge Management Center Illuminating the Druggable Genome Knowledge Management Center (IDG KMC) provides an integrated informatics solution encompassing accrual, storage, analysis, visualization and dissemination of data and knowledge related to unannotated human proteins, focused on G-protein-coupled receptors, kinases, ion channels and nuclear receptors. IDG KMC offers informatics resources to assist investigators in uncovering the biological functions of previously uncharacterized proteins, by identifying knowledge gaps and prioritizing candidates for further experiment using multiple technology platforms, as outlined at the IDG website, http://targetcentral.ws/. These resources are disseminated via the IDG KMC portal, Pharos, https://pharos.nih.gov/ which supports efficient query, browsing and analysis of the IDG KMC

3. Text mining resources from JensenLab (CPR): http://jensenlab.org/ 4. Gene Reference into Function (NCBI): http://www.ncbi.nlm.nih.gov/gene/about-generif 5. Antibodypedia website: http://www.antibodypedia.com/

The potency cut-offs for small-molecule modulators used for the target classes are: 9. Rouillard, A.D. et al. The harmonizome: a collection of processed datasets gathered to serve and mine knowledge about genes and proteins. Database http://dx.doi.org/10.1093/database/baw100 (2016)

University of Miami, 1600 NW 10th Avenue, RMSB 6112B, Miami, Florida 33136, USA.

12. ChEMBL: https://www.ebi.ac.uk/chembl/