

# Thera-SAbDab: the Therapeutic Structural Antibody Database

Matthew I.J. Raybould<sup>1</sup>, Claire Marks<sup>1</sup>, Alan P. Lewis<sup>2</sup>, Jiye Shi<sup>3</sup>, Alexander Bujotzek<sup>4</sup>, Bruck Taddese<sup>5</sup> and Charlotte M. Deane<sup>1,\*</sup>

<sup>1</sup>Oxford Protein Informatics Group, Department of Statistics, University of Oxford, 24-29 St Giles', Oxford OX1 3LB, UK, <sup>2</sup>Data and Computational Sciences, GlaxoSmithKline Research and Development, Gunnels Wood Road, Stevenage SG1 2NY, UK, <sup>3</sup>Chemistry Department, UCB Pharma, 216 Bath Road, Slough SL1 3WE, UK, <sup>4</sup>Roche Pharma Research and Early Development, Large Molecule Research, Roche Innovation Center Munich, DE-82377 Penzberg, Germany and <sup>5</sup>Discovery Sciences Department, AstraZeneca, Granta Park, Cambridge CB21 6GH, UK

Received August 08, 2019; Revised September 09, 2019; Editorial Decision September 13, 2019; Accepted September 24, 2019

## ABSTRACT

The Therapeutic Structural Antibody Database (Thera-SAbDab; <http://opig.stats.ox.ac.uk/webapps/therasabdab>) tracks all antibody- and nanobody-related therapeutics recognized by the World Health Organisation (WHO), and identifies any corresponding structures in the Structural Antibody Database (SAbDab) with near-exact or exact variable domain sequence matches. Thera-SAbDab is synchronized with SAbDab to update weekly, reflecting new Protein Data Bank entries and the availability of new sequence data published by the WHO. Each therapeutic summary page lists structural coverage (with links to the appropriate SAbDab entries), alignments showing where any near-matches deviate in sequence, and accompanying metadata, such as intended target and investigated conditions. Thera-SAbDab can be queried by therapeutic name, by a combination of metadata, or by variable domain sequence - returning all therapeutics that are within a specified sequence identity over a specified region of the query. The sequences of all therapeutics listed in Thera-SAbDab (461 unique molecules, as of 5 August 2019) are downloadable as a single file with accompanying metadata.

## INTRODUCTION

Immunotherapeutics derived from B-cell genes are an increasingly successful and significant proportion of the global drugs market, designed to treat a wide range of diseases (1–3).

Whole monoclonal antibody (mAb) therapies dominate the industry - drugs that mimic natural antibodies by containing two identical variable domain structures with a particular specificity (3). The broader class of monoclonal therapies also includes Fragment antigen binding (Fab) regions (a single arm of a whole antibody), single-chain Fv (scFv) regions (a heavy and light chain variable domain connected by an engineered glycine-rich linker), and single-domain variable fragments. These fragments can be expressed in dimeric form to improve avidity, or conjugated with polyethylene glycol ('pegylated') for slower clearance (4), with radioisotopes for diagnostic purposes (5), or with radioisotopes or noxious small molecules/peptides for cytotoxicity (6).

Recent developments in protein engineering have resulted in bispecific immunotherapies, where two distinct variable domain binding sites are incorporated into a single protein. As of June 2019, bispecific mAbs, linked Fabs, linked scFvs and linked single-domain variable fragments have all been assessed in clinical trials (7).

A primary source of information on immunotherapies is the World Health Organisation (WHO), which publishes biannual 'Proposed' (8) and 'Recommended' (9) International Nonproprietary Name (INN) lists. These INNs serve as globally-recognized generic names by which pharmaceuticals can be identified. To be granted an INN, applicants must include a full amino acid sequence, the closest V and J gene, the IG subclass, and the light chain type (see [https://extranet.who.int/tools/inn\\_online\\_application/](https://extranet.who.int/tools/inn_online_application/)). This information, coupled with the \$12 000 cost of application (as of August 2019), makes INN lists a useful source of therapies that companies intend to carry forward into clinical trials.

Several databases already harvest this information. Two non-commercial antibody-specific resources are the IMGT

\*To whom correspondence should be addressed. Tel: +44 1865 272860; Email: [deane@stats.ox.ac.uk](mailto:deane@stats.ox.ac.uk)

Monoclonal Antibody Database (IMGT mAb-DB; <http://www.imgt.org/mAb-DB> (10), and WHOINNIG (<http://www.bioinf.org.uk/abs/abybank/whoinnig>).

The Therapeutic Antibody Database (TABS; <https://tabs.craic.com>) is antibody-specific and commercial, also scraping patents for therapies. Other databases not specific to antibodies can also capture WHO information, such as ChEMBL (<https://www.ebi.ac.uk/chembl>), DrugBank (<https://www.drugbank.ca>) and KEGG DRUG (<https://www.genome.jp/kegg/drug>).

Most databases supply additional metadata for their therapeutic entries, such as clinical trial status, companies involved in development, target specificity, and alternative names. For example, the recently published ABCD database provides antibody synonyms, antigen UniProt links and publication references (11). However, while these repositories supply sequence information (either on individual summary pages or through reference to the primary literature), it is currently not possible to query them by sequence, nor to bulk-download relevant sets of therapeutic sequences for direct bioinformatic analysis.

Structural knowledge about both the intended target and the therapeutic lead compound is of high importance for rational drug discovery (12,13). For example, co-crystal complexes reveal where a drug binds to its target (the surface ‘epitope’), and separately-solved structures enable more accurate docking experiments. It can also assist subsequent development and optimization, as homology models of mutants derived from a known structure are in general more accurate than those for which no close structural partner is available (14). The Protein Data Bank (15) (PDB) now contains over 150 000 solved structures, and though it is highly biased towards certain protein classes, many diverse targets of pharmacological interest are represented. A significant fraction of these structures contain antibody variable domains, and these are recorded by the Structural Antibody Database (SAbDab (16); 7184 variable domain structures over 3663 PDB entries as of 5 August 2019). Both IMGT mAb-DB and TABS report a set of known therapeutic structures in the PDB, but their reported structural coverage of therapeutic space is low. For example, neither database reports any known structural information for bispecific immunotherapeutics.

To address these deficiencies, we have created the Therapeutic Structural Antibody Database (Thera-SAbDab; <http://opig.stats.ox.ac.uk/webapps/therasabdab>). We harvest sequences as they are released by the WHO, number them with ANARCI (17), and perform a weekly sequence alignment of all therapeutic variable domain sequences to the sequences of known structures stored in SAbDab. Structures with sequence identity matches of 100%, 99% and 95–98% are recorded and categorized, with alignments on each therapeutic summary page to show precisely where each near-identical structure differs from the therapeutic sequence.

Thera-SAbDab can be queried by INN, by a combination of metadata, such as INN proposal year, clinical trial status, or target, or by sequence (including over a specified region of the sequence). We make available all therapeutic sequences contained within Thera-SAbDab, alongside metadata, to facilitate further research.

## DATA SOURCES

### Sequence data

Proposed INN lists (8,9), published by the WHO, are the source of the majority of sequence information in Thera-SAbDab. These are released biannually (one in January/February and another in June/July) and—since list P95 in 2006—represent a reliable record of variable domain sequences for all antibody- and nanobody-related therapeutics granted a proposed INN. Of the 129 antibody-related therapeutics proposed before 2006, we were able to find sequence information for 47 (36.4%) through the IMGT mAb-DB (<http://www.imgt.org/mAb-DB/>). Although we continue to search, and joint academia-industry initiatives such as Abvance encourage their release (<https://www.pistoiaalliance.org/projects/abvance/>), sequences for the remaining 82 may never become public knowledge.

All sequences are then numbered by ANARCI (17), which uses Hidden Markov Models to align input sequences to pre-numbered germline sequences. Assigning a numbering allows users to more easily interpret the significance of mutations in near-identical sequence matches. For example, if the mismatch occurs in the extremities of the framework region, it may be judged to have minimal effect on binding site structure.

### Structural data

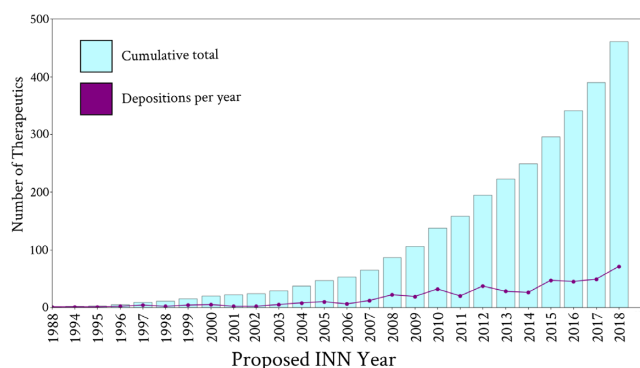
Thera-SAbDab compares all numbered therapeutic sequences to the structures in SAbDab (16), which prefilters the PDB (15) for all structures whose sequences align to B-cell germline genes. As all SAbDab structures are also pre-numbered, the comparison of therapeutics to public structural space is efficient. All the existing functionality of SAbDab (e.g. interactive molecular viewers and numbered structure downloads) is made easily accessible from Thera-SAbDab search results.

### Therapeutic metadata

Therapeutic metadata comprises a mixture of inherent characteristics and continually-changing status updates.

Certain static properties can be acquired automatically. For example, light chain type is identified through our ANARCI germline alignment (17), while isotype, INN Proposed and Recommended years, and intended target(s) can be harvested directly from the INN lists. Sequence comparison can also be used to identify where different INN names refer to identical variable domains. Other characteristics, such as which companies are involved in therapeutic development, must be manually curated at the time of deposition.

Time-dependent characteristics for new entries are also manually curated after sequence identification, and thereafter every 3 months. We source clinical trial information, developmental status, and investigated condition data from a range of sources including AdisInsight (<https://adisinsight.springer.com>), ClinicalTrials.gov (<https://clinicaltrials.gov>), and DrugBank (<https://www.drugbank.ca>). These websites are updated more regularly, and so are preferable sources for this time-sensitive metadata; we



**Figure 1.** The number of antibody- and nanobody-related therapeutics assigned an International Nonproprietary Name (INN) by year. A record number of 72 of these therapeutics were recognized by the WHO in 2018.

include these fields in Thera-SAbDab to allow for more pharmacologically-relevant searches, as well as to identify all post Phase-I candidates for inclusion in our five updating developability guidelines (18).

## CONTENTS

As of 5 August 2019, Thera-SAbDab is tracking 558 INNs, representing 543 unique therapeutics. Of the 558 INN names, 473 could be mapped to variable domain sequences (87.1%), representing 461 unique therapeutics with sequence data. 436 were monoclonal therapies (three pairs of which share identical variable domains: avelumab & bintrafusp, losatuxizumab & serclutamab and radretumab & bifikafusp), and 25 were bispecific therapies. Plotting the cumulative sum of these unique therapeutics by year deposited in a WHO ‘Proposed INN’ list shows an exponential increase since the early 2000s (Figure 1).

We searched the IMGT mAb-DB (10) and TABS databases (on 28 June 2019) for structures of these 461 therapeutics. IMGT mAb-DB identified 72 structures of therapeutic variable domains, across 36 different monoclonal therapeutics, while TABS reported 53 structures of therapeutic variable domains, across 32 different monoclonal therapeutics. In contrast, Thera-SAbDab (at the 100% sequence identical threshold) contained 152 therapeutic variable domain structures, across 84 distinct monoclonal therapeutics and 7 distinct bispecific therapeutics. A further 21 monoclonal therapeutics had maximum sequence identity matches of 99% (up to two mutations away from a publicly-available structure), and 13 monoclonals and 4 bispecifics had maximum sequence identity matches of 95–98%. We conclude that, at present, around a quarter (27.1%) of WHO-recognized monoclonal therapeutics have exact or close ( $\geq 95\%$  sequence identity) structural coverage. 44.0% of bispecific therapeutics have at least one variable domain with exact or close structural coverage, and two have exact matches for both variable domains.

Thera-SAbDab contains structural information for even the most diversely-formatted therapeutics. Ozoralizumab, a bispecific therapy in active Phase-III clinical trials for rheumatoid arthritis, has a VH(TNFA)–VH(ALB)–VH(TNFA) configuration, where VH(TNFA) is a heavy

chain designed to bind to TNF- $\alpha$ , and VH(ALB) is another heavy chain designed to bind ALB. Thera-SAbDab has identified a structure for the TNFA binding domain with sequence identity of 95.65% [5m2j; chain D]. Inspection of the sequence alignment shows that 5m2j has a 100% Chothia-defined CDRH3 sequence match to VH(TNFA), and in fact only differs by one mutation across all Chothia-defined (19) CDRs: 31D in VH(TNFA) is 31N in 5m2j. 5m2j is a VHH2 llama nanobody, suggesting that SAbDab’s coverage of nanobody structural space will be increasingly highlighted by Thera-SAbDab as more single-chain therapies arrive in the clinic.

Therapeutically-relevant structures are continually being deposited in the PDB, even many years after initial development. For example, since 2009, the WHO have recorded nine antibody-related therapeutics against IL17A—seven monoclonals and two bispecifics. The first, secukinumab, was recognized in 2009, and since 2014 has been approved for use in certain types of arthritis, psoriasis, and spondylitis. As of early June 2019, there were no close structures for any of these IL17A-binders. However, on 19 June 2019, Eli Lilly deposited an exact variable domain structure for ixekizumab (an IL17A-targeting monoclonal antibody, 6nov) and a close structure for tibilizumab (an IL17A-binding and TNFSF13B-binding bispecific antibody, 6nou) in the PDB (20). SAbDab detected and numbered them in its weekly update, making Thera-SAbDab the first antibody database to link to the structures of IL17A-binding therapeutic antibodies.

## USAGE

There are multiple ways to search Thera-SAbDab. Thera-SAbDab can be queried directly by INN if structural information about a particular therapeutic is needed. Alternatively a combination of metadata can be specified to identify structures for a particular subset of therapeutic space, for example binders to a particular antigen, or therapeutics at a particular stage of clinical trials (Figure 2A). Results are returned in a table format, with links to each therapeutic summary page and a selected array of metadata (Figure 2B).

Each therapeutic summary page lists a structural summary (including our database sequence), with links to relevant SAbDab entries (with PDB codes and chains), and alignment charts (if structures with 95–99% sequence identity are detected). Each SAbDab link redirects the user to the SAbDab summary page for the relevant PDB entry, where all existing functionality can be accessed. Links to appropriate SAbPred (21) informatics tools (such as ABody-Builder (22) for variable domain structure modelling, and TAP (18) for developability assessment) are also provided. Finally, we list all the remaining metadata that we have recorded for the therapeutic, ranging from records of investigated conditions, to which companies are developing the therapeutic, to its estimated developmental status.

A third way to search Thera-SAbDab is by sequence (Figure 2C and D). This can be harnessed in numerous ways. For example, by querying with a known therapeutic sequence, researchers can look for sequence commonalities between therapeutics over any region of the variable domain. Alter-



**Figure 2.** Searching Thera-SABDab. (A) Search by attribute. Here, we search for any therapeutic designed to bind to ERBB2 (often over-expressed in breast cancer). (B) Eight therapeutics are designed to bind to ERBB2, seven monoclonals and one bispecific. Four have exact structural information for the ERBB2 binding site. Click the therapeutic name to enter the therapeutic summary page. (C) Search by sequence. Here we search for therapeutics with at least 70% sequence identity across the heavy and light chain CDRs of the input sequence. (D) Any results are returned alongside sequence identity across the specified region. Alignments show any sequence mismatches across the variable domain sequence.



natively, by querying with a developmental candidate sequence, researchers can search for similarity to any other therapeutic, or specifically to those designed to bind to the same target. This could identify potential patenting issues, highlight a risk of polyspecificity, or suggest a binding mode to the intended target.

A further selection of sample use cases for Thera-SAbDab are available at <http://opig.stats.ox.ac.uk/webapps/therasabdab/about>.

## ACCESSIBILITY OF THE DATA

Thera-SAbDab can be queried at <http://opig.stats.ox.ac.uk/webapps/therasabdab>. All sequence data harvested by Thera-SAbDab can be downloaded from the 'Downloads' tab of the search page. Sequences are supplied alongside the therapeutic INN, format, isotype, light chain category, highest clinical trial stage reached, and estimated developmental status. We also supply a list of therapeutics for which sequence information has not yet been released.

## CONCLUSION

We have created Thera-SAbDab with the central aim of collating all public structural knowledge for WHO-recognized antibody- and nanobody-related therapeutic variable domains. Rather than relying on text-mining approaches, which can miss PDB depositions that omit reference to the structure's therapeutic relevance, Thera-SAbDab uses a systematic approach at the level of sequence identity to detect exact and close matches to our repository of therapeutic variable domains.

This approach has not only enabled us to identify over twice the number of monoclonal therapies with 100% sequence-identical structures in the PDB than in existing databases, but has also identified exact variable domain structures for several bispecific therapies. Our approach can also distinguish between PDB structures with 100%, 99%, and 95–98% sequence identity matches. Sequence alignments guide the interpretation of structures of near-identical sequence.

Like IMGT-DB, Thera-SAbDab can be queried by metadata, but uniquely it can also be queried by variable domain sequence. This enables researchers to identify any therapeutics proximal over any variable domain region to their query sequence.

Thera-SAbDab's sequence database will be updated with new sequence information twice per year, in line with the release of new WHO Proposed INN lists. An updated list of all therapeutic variable domain sequences with metadata is supplied as a single file to facilitate further analysis, for example into the properties of therapeutic antibody-antigen interfaces.

As shown for IL17A-binding therapeutics, new clinically-relevant structures are continually being released. Accordingly, Thera-SAbDab checks SAbDab after each weekly update for new matches, ensuring that this data is rapidly captured.

## FUNDING

Engineering and Physical Sciences Research Council and Medical Research Council [EP/L016044/1]; Glaxo-SmithKline plc; AstraZeneca plc; F. Hoffmann-La Roche AG; UCB Celltech. Funding for open access charge: RCUK Open Access Block Grant and processed through the Bodleian Library, University of Oxford.

*Conflict of interest statement.* None declared.

## REFERENCES

- Grilo, A.L. and Mantalaris, A. (2018) The increasingly human and profitable monoclonal antibody market. *Trends Biotechnol.*, **37**, 9–16.
- Steele, S., Vandenbroucke, R.E. and Libert, C. (2017) Nanobodies as therapeutics: big opportunities for small antibodies. *Drug Discov. Today*, **21**, 1076–1113.
- Kaplon, H. and Reichert, J.M. (2019) Antibodies to watch in 2019. *mAbs*, **11**, 219–238.
- Jevšvar, S., Kusterle, M. and Kenig, M. (2012) PEGylation of Antibody Fragments for Half-Life Extension. In: Proetzel, G and Ebersbach, H (eds). *Antibody Methods and Protocols. Methods in Molecular Biology (Methods and Protocols)*. Vol. **901**, Humana Press, Totowa.
- Steiner, M. and Neri, D. (2011) Antibody-radionuclide conjugates for cancer therapy: historical considerations and new trends. *Clin. Cancer Res.*, **17**, 6406–6416.
- Beck, A., Goetsch, L., Dumontet, C. and Corvaia, N. (2017) Strategies and challenges for the next generation of antibody-drug conjugates. *Nat. Rev. Drug Discov.*, **16**, 315–337.
- Labrijn, A.F., Janmaat, M.L., Reichert, J.M. and Parren, P.W.H.I. (2019) Bispecific antibodies: a mechanistic review of the pipeline. *Nat. Rev. Drug Discov.*, **18**, 585–608.
- WHO (2018) Proposed International Nonproprietary Names (INN) List 120. *WHO Drug Information*, **32**, 559–689.
- WHO (2019) Recommended International Nonproprietary Names (INN) List 81. *WHO Drug Information*, **33**, 59–134.
- Poirion, C., Wu, Y., Ginestoux, C., Ehrenmann, F., Duroux, P. and Lefranc, M.-P. (2010) IMGT/mAb-DB: the IMGT database for therapeutic monoclonal antibodies. *JOIM* **2010**, **13**, 470b.
- Lima, W. C., Gasteiger, E., Marcantili, P., Duek, P., Bairoch, A. and Cosson, P. (2019) The ABCD database: a repository for chemically defined antibodies. *Nucleic Acids Res.*, **48**, gkz714.
- van Montfort, R.L.M. and Workman, P. (2017) Structure-based drug design: aiming for a perfect fit. *Essays Biochem.*, **61**, 431–437.
- Raybould, M.I.J., Wong, W.K. and Deane, C.M. (2019) Antibody-antigen complex modelling in the era of immunoglobulin repertoire sequencing. *Mol. Syst. Des. Eng.*, **4**, 679–688.
- Muhammed, M.T. and Aki-Yalcin, E. (2019) Homology modeling in drug discovery: overview, current applications, and future perspectives. *Chem. Biol. Drug Des.*, **93**, 12–20.
- Berman, H.B., Westbrook, J., Feng, Z., Gilliland, G., Bhat, T.N., Weissig, H., Shindyalov, I.N. and Bourne, P.E. (2000) The Protein Data Bank. *Nucleic Acids Res.*, **28**, 235–242.
- Dunbar, J., Krawczyk, K., Leem, J., Baker, T., Fuchs, A., Georges, G., Shi, J. and Deane, C.M. (2014) SAbDab: the structural antibody database. *Nucleic Acids Res.*, **42**, D1140–D1146.
- Dunbar, J. and Deane, C.M. (2016) ANARCI: antigen receptor numbering and receptor classification. *Bioinformatics*, **32**, 298–300.
- Raybould, M.I.J., Marks, C., Krawczyk, K., Taddese, B., Nowak, J., Lewis, A.P., Bujotzek, A., Shi, J. and Deane, C.M. (2019) Five computational developability guidelines for therapeutic antibody profiling. *Proc. Natl. Acad. Sci. U.S.A.*, **116**, 4025–4030.
- Al-Lazikani, B., Lesk, A.M. and Chothia, C. (1997) Standard conformations for the canonical structures of immunoglobulins. *J. Mol. Biol.*, **273**, 927–948.
- Benschop, R.J., Chow, C.-K., Tian, Y., Nelson, J., Barmettler, B., Atwell, S., Clawson, D., Chai, Q., Jones, B., Fitchett, J. et al. (2019) Development of tibatuzumab, a tetravalent bispecific antibody targeting BAFF and IL-17A for the treatment of autoimmune disease. *mAbs*, **11**, 1175–1190.

21. Dunbar,J., Krawczyk,K., Leem,J., Marks,C., Nowak,J., Regep,C., Georges,G., Kelm,S., Popovic,B. and Deane,C.M. (2014) SAbPred: a structure-based antibody prediction server. *Nucleic Acids Res.*, **44**, W474–W478.
22. Leem,J., Dunbar,J., Georges,G., Shi,J. and Deane,C.M. (2016) ABodyBuilder: automated antibody structure prediction with data-driven accuracy estimation. *mAbs*, **8**, 1259–1268.