



Extending the small-molecule similarity principle to all levels of biology with the Chemical Checker

Miquel Duran-Frigola^{1,3}✉, Eduardo Pauls^{1,3}, Oriol Guitart-Pla¹, Martino Bertonì¹, Víctor Alcalde¹, David Amat¹, Teresa Juan-Blanco¹ and Patrick Aloy^{1,2}✉

Small molecules are usually compared by their chemical structure, but there is no unified analytic framework for representing and comparing their biological activity. We present the Chemical Checker (CC), which provides processed, harmonized and integrated bioactivity data on ~800,000 small molecules. The CC divides data into five levels of increasing complexity, from the chemical properties of compounds to their clinical outcomes. In between, it includes targets, off-targets, networks and cell-level information, such as omics data, growth inhibition and morphology. Bioactivity data are expressed in a vector format, extending the concept of chemical similarity to similarity between bioactivity signatures. We show how CC signatures can aid drug discovery tasks, including target identification and library characterization. We also demonstrate the discovery of compounds that reverse and mimic biological signatures of disease models and genetic perturbations in cases that could not be addressed using chemical information alone. Overall, the CC signatures facilitate the conversion of bioactivity data to a format that is readily amenable to machine learning methods.

The current catalog of purchasable chemical substances amounts to 100 million¹, and databases containing bioactivity data annotate a few million of them^{2,3}. The deluge of publicly available data resources has transformed the field of pharmacology⁴, although, with the exception of metabolomics, omics-based biomedical research continues to be acutely gene-centric and difficult to link to chemical compounds⁵. The limited availability of systematic, comprehensive small-molecule datasets greatly hampers the discovery of compound–biomolecule interactions and their possible links to disease. In consequence, often the first step in characterizing the bioactivity of a compound is to look for structurally similar molecules^{5,6}. The so-called ‘similarity principle’ has become the driving force of drug discovery: most known drugs were inspired by natural products^{7,8}; chemical libraries are created by combining or decorating privileged chemotypes⁹ and the design of lead drug candidates departs from hit compounds identified in experimental screening assays¹⁰. Thus, compound similarities are the primary measure to chart and exploit chemical space.

The release of compound databases has led to the realization that the similarity principle applies beyond chemical properties. For instance, molecules with similar cell-sensitivity profiles tend to share the mechanism of action^{11,12}, as do drugs eliciting similar side effects¹³, even when their chemical structures are unrelated. Hence, biological similarities offer an alternative means of functionally characterizing small molecules, potentially to a degree that is closer to clinical observations and beyond the mere inspection of chemical analogs¹⁴. However, there is no convention to compare the biological profiles of small molecules, since available bioactivity data are sparse, incomplete and often of dubious quality¹⁵, requiring thorough preprocessing and integration. As a result, the extent to which the similarity principle can be generalized to biology (and possibly embrace omics techniques) remains unclear. In this article, we present the CC, a resource that expands the similarity principle along the drug discovery pipeline from in vitro assays to clinical

observations by treating bioactivity data within a unified analytical framework. To illustrate the capabilities of the CC in the day-to-day drug discovery process, we describe applications that include chemical library visualization, identification of compounds reverting disease-associated signatures and discovery of small molecules that mimic the biological effect of approved biologics.

Results

Five levels of complexity for small-molecule data. Of all existing compounds, approved drugs (APD) are probably the most widely characterized¹⁶. Small-molecule data can be organized in five levels of increasing complexity, based on the principal steps of the drug discovery process (Fig. 1a). A drug is often an organic molecule (A, chemistry) that interacts with one or several protein receptors (B, targets), triggering perturbations of biological pathways (C, networks) and eliciting phenotypic outcomes that can be measured in, for example, cell-based assays (D, cells) before delivery to patients (E, clinics). We used these five categories to classify the information stored in chief compound databases, including chemogenomics resources, cell-based screens and, when available, clinical reports of drug effects (Methods).

We then divided each level (A–E) into five sublevels (1–5) corresponding to distinct types or scopes of the data. In total, the CC contains 25 well-defined categories meant to illustrate the most relevant aspects of small-molecule characterization. In particular, we stored the two-dimensional (2D) (A1) and three-dimensional (3D) (A2) structures of compounds, together with their scaffolds (A3), functional groups (structural keys) (A4) and physicochemical parameters (A5). We also retrieved therapeutic targets (mechanisms of action) (B1) and drug-metabolizing enzymes (B2), and molecules cocrystallized with protein chains (B3). We incorporated literature binding data (B4) from chief chemogenomics databases, and high-throughput target screening results (B5). Moving to a higher order of biology, we looked for ontological classifications

¹Joint IRB-BSC-CRG Program in Computational Biology, Institute for Research in Biomedicine (IRB Barcelona), The Barcelona Institute of Science and Technology, Barcelona, Spain. ²Institució Catalana de Recerca i Estudis Avançats (ICREA), Barcelona, Spain. ³These authors contributed equally: Miquel Duran-Frigola, Eduardo Pauls. ✉e-mail: miquel.duran@irbbarcelona.org; patrick.aloy@irbbarcelona.org

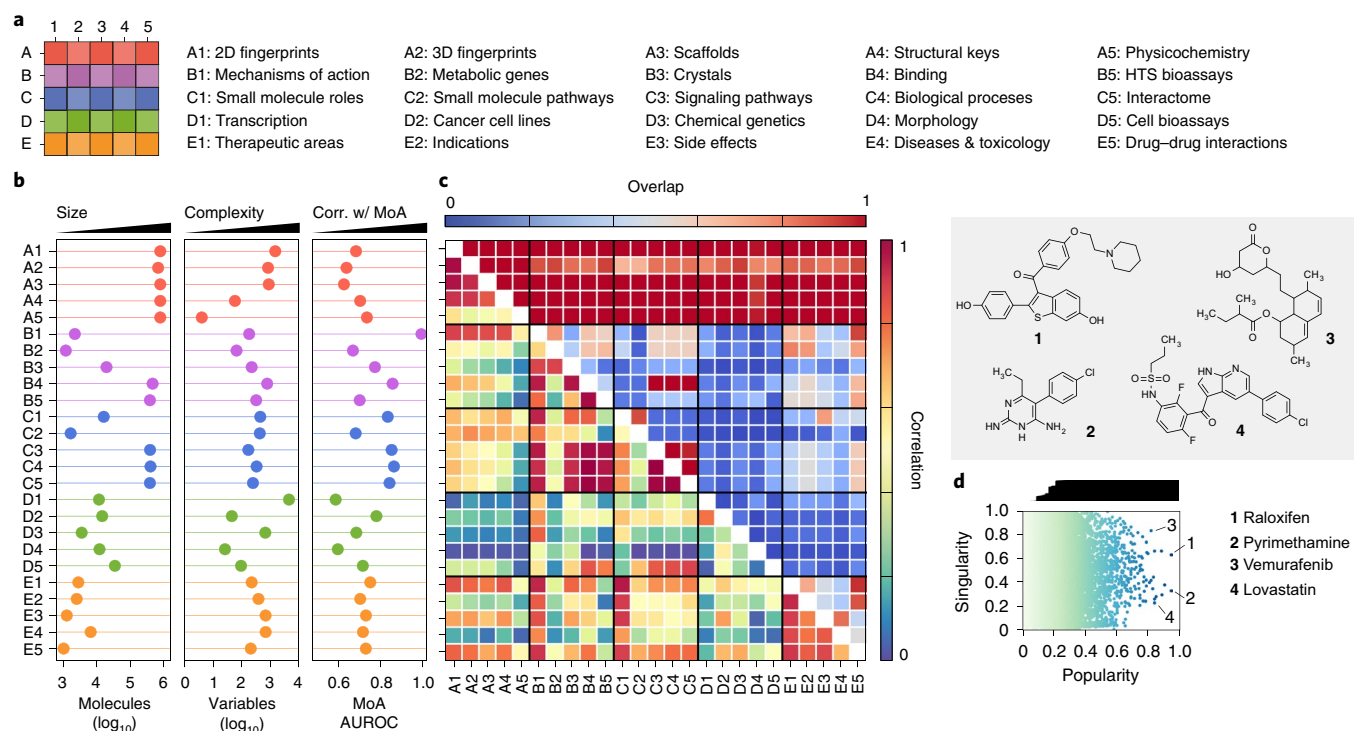


Fig. 1 | CC statistics. **a**, The organization of the 5 × 5 CC spaces. **b**, Number of molecules (size), signature length (that is, number of latent variables as a measure of data complexity) and area under the receiver operating characteristic (AUROC) curve performances when checking whether similar molecules in each CC space tend to share mechanism of action (MoA). **c**, Overlap between CC spaces, in terms of number of shared molecules (upper triangle) and correlation k between CC spaces (lower triangle). **d**, Popularity and singularity of molecules. Popularity refers to the proportion of CC spaces in which the molecule is present (correcting for correlation between CC spaces), and singularity refers to the ‘uniqueness’ of the molecule. The larger the number of molecules showing similarity to a given molecule, the less singular the molecule is. Popular molecules within a wide range of singularities are highlighted. For example, raloxifen (**1**), pyrimethamine (**2**) and vemurafenib (**3**) have data in many CC spaces. Likewise, some molecules are more singular than others for which many analogs exist throughout the CC organization (for example, lovastatin (**4**)). HTS, high-throughput screening; Corr., correlation.

of compounds (small molecule roles) (C1) and focused on human metabolites in a genome-scale metabolic network (small molecule pathways) (C2). In addition, we kept the signaling pathways (C3), biological processes (C4) and protein–protein interactions (PPIs) (C5) of the previously collected binding data. To capture cell-level information, we gathered differential gene expression profiles (transcriptomics) (D1) and compound growth-inhibition potencies across cancer cell lines (D2). Similarly, we gathered sensitivity profiles over an array of yeast mutants (chemical genetics) (D3), as well as cell morphology changes (high-content screening) (D4). Additional cell-sensitivity data available from the literature were also collected (D5). To organize clinical data, we used the traditional Anatomical Therapeutic Chemical (ATC) classification of drugs (therapeutic areas) (E1), and also drug indications (E2) and side effects (E3) expressed as disease terms, together with therapeutic or adverse outcomes of molecules other than drugs such as environmental chemicals (disease and toxicology) (E4). Finally, we included drug–drug interactions (DDIs) known to raise pharmacokinetic and efficacy issues (E5).

Further rationale for the choice of the 25 CC categories is presented in Table 1. Overall, we believe that the CC organization is a good representation of what is known of small molecules in the public domain (Supplementary Table 1). In the Methods, we extensively describe the data collection protocol. We adopted well-accepted standards, harmonized chemical entries and filtered bioactivities (Supplementary Fig. 1). For example, in the CC D1 space, we discarded those molecules whose transcriptional response was not noticeable and, similarly, only notorious distortions of cell morphology were kept in D4, excluding innocuous compounds.

Likewise, we applied target-class-specific potency cutoffs to binding data¹⁷. At the ‘networks’ level (C), we incorporated ontologies and systems biology datasets that are typically outside the scope of compound databases.

The CC contains and catalogs information on nearly 800,000 bioactive compounds (Fig. 1b), and is mainly focused on human pharmacology data (Supplementary Fig. 2). Evidently, fewer molecules are available as we advance along the CC levels from A to E: **chemical information (A) is always available (778,460 molecules), whereas clinical data (E) are scarce (9,165 molecules, including 4,232 drugs)**. Most molecules come from the binding literature (B4) or target-based high-throughput screening bioassays (B5) (adding up to 705,685 entries in B), and part of this knowledge is transferred to network levels (C3–5) by virtue of biological ontologies, pathways and PPIs. On a similar scale, the current throughput of cell-based assays (D1–4) is of about 10,000–20,000 molecules.

Signature-based representation of the data. **Inspired by the success of chemical descriptors and fingerprints to represent compound structures¹⁸**, we chose to express bioactivity data in a common vector format. Details on how we obtained vectors (signatures) for the 25 CC spaces are given in the Methods. **In brief, we treated categorical data as sets of terms**, these being proteins, pathways, ATC codes, bit positions of a chemical fingerprint and so on. We then removed frequent and rare terms, and down-weighted the less informative ones (that is, promiscuous targets, generic biological processes and so on). Finally, we applied a dimensionality reduction technique, called latent semantic indexing (LSI), to ensure that signature components were orthogonal and sorted by their

Table 1 | Rationale for the choice of the five CC levels (A–E) and sublevels (1–5). Explanations are given from the perspective of drug discovery

| Level | Name | Rationale |
|-------|-----------|---|
| A | Chemistry | Database search and large-scale prediction tools typically use 2D encodings of compounds (A1). Target prediction algorithms often require 3D representations of the compounds (A2), which usually involves an energy-optimization step. In addition, a convenient way to browse the chemical space is through the inspection of scaffolds (A3), and a means to communicate with synthetic chemists is through structural keys (functional groups) (A4). Finally, physicochemical parameters such as the molecular weight are used to rapidly characterize small-molecule entities, together with drug-likeness estimations (A5). |
| B | Targets | For a relatively small number of molecules (drugs), targets with pharmacological action are known (B1) and drug-metabolizing enzymes, transporters, and carriers (B2) are crucial determinants of drug safety (and efficacy). Another prominent set of small molecules are those that have been cocrystallized with protein chains (B3), as they greatly inform of structure-based molecular design. Beyond these, there is a large corpus of binding affinity measurements (B4) and target functional assays (B5) available from the literature and screening campaigns. |
| C | Networks | Biologists put molecules in context via pathways, ontologies and networks. The biological roles of eminent chemical entities are part of an ontology (C1). Some entities that are metabolites can be found in the human metabolic network (C2), where substrates and products of enzymes are linked by reactions. However, these 'higher-order' annotations correspond to a minority of molecules. To include systems-level data for more compounds, one can incorporate large-scale compound-protein interaction data (for example, B4), and the canonical pathways (C3) and biological processes (C4) of the proteins may be kept, correspondingly, as annotations for the compound (guilt-by-association). With finer detail, the neighborhoods of these proteins in PPI networks can be inspected as well (C5). |
| D | Cells | Cell-based assays are bringing about the largest increase in the variety of relevant data. The LINCS consortium collects gene expression changes after compound dosage (D1), and pioneering approaches, such as the NCI-60, continue to produce sensitivity profiles of cancer cell line panels (D2). Similarly, chemical genetics profiles have been proposed to complement genetic interactions in yeast (D3). A very different type of experiment is high-content phenotypic screening, in which morphological changes of cells (D4) are measured with microscopy. Growth and proliferation assays (D5) accumulated in the literature over the years can be added to these techniques. |
| E | Clinics | Clinical data represent the last level of complexity. Drug molecules have been traditionally classified using a hierarchical taxonomy based on anatomy and therapeutic areas (E1), although medical vocabularies may be adopted to better defining drug indications (E2) and linking them to disease genetics studies. Likewise, side effects (E3) can be cataloged by parsing drug package inserts and, with a broader scope, there are resources that mine the literature for beneficial and harmful associations between compounds and diseases, including molecules other than drugs, such as environmental chemicals (E4). Finally, acknowledging DDIs (E5) is crucial for medical prescription and the avoidance of undesired pharmacokinetics and toxicity. |

contribution to explaining the 'variance' of the data. An analogous procedure (that is, robust scaling followed by a principal component analysis (PCA)) was performed on continuous data. For each CC space, we kept the number of components retaining 90% of the variance (Supplementary Fig. 3a). **As a result, we obtained 25 numerical matrices, rows corresponding to molecules and columns composing signatures. We named these CC vectors 'type I signatures.'**

Most type I signatures have a length between 500 and 1,500 components (Fig. 1b). Longer signatures denote higher complexity or sparseness of the data. Signatures based on gene expression (D1) are the longest, followed by binding data (B4 and B5) and the fine-grained chemical descriptors (A1 and A2). Conversely, physicochemical (A5) and cancer cell line sensitivity signatures (D2) are the shortest. Of note, morphology (D4) signatures require only 26 components to account for the original 812 features, indicating high interdependency of raw measurements (Supplementary Fig. 3b).

We observed that, in all 25 CC spaces, compounds with similar signatures tend to share the mechanism of action and therapeutic area (Fig. 1b and Supplementary Fig. 4a). Indeed, bioactivity signatures often correlate better with known mechanisms of action than the more classical chemical signatures (A1–5). Pairwise similarity measurements reveal clusters of molecules and molecules in the same cluster share targets and therapeutic areas (Supplementary Fig. 4b). More generally, using a similarity-based correlation analysis (Methods) we certified an inter-connection between CC spaces (Fig. 1c and Supplementary Fig. 5). Certain links within the CC were expected by design, such as connections within the chemistry spaces (A1–4), or those between binding data (B4) and functionally related versions of them (C3–5). Other correlations have a straightforward interpretation (for example, drugs with similar

targets (B1) have similar indications (E1–2)), and some reflect recognized research biases. For example, we found stronger links between chemistry and mechanisms of action (B1) and therapeutic indications (E1), compared to spaces representing collateral processes such as metabolic enzyme interactions (B2) and toxicology events (E4). Notably, we observed correlations between unbiased (omics) datasets. For instance, we found cell-sensitivity profiles (D2) to be linked to many CC levels, including the clinical (E) ones (even though this connection is of particular relevance, identifying the main drivers for it is outside the scope of the current study, and would require further analysis and validation). In turn, D2 appeared to be complementary to a comparable yeast sensitivity screening panel (D3) (Supplementary Fig. 6), suggesting that incorporating cross-species data could further enrich the CC cellular (D) layers.

To balance the numerical complexity across CC spaces, we derived an embedded (128-dimension) version of CC signatures (type II signatures). This was achieved by first building similarity networks based on type I signatures, **and then using a network embedding technique to capture (embed) the vicinity of each node (molecule) in a vector space, so that local similarities and more global network properties are seized** (Methods, Fig. 2a and Supplementary Fig. 7). Figure 2b displays type II signatures for five representative CC datasets, related to drugs used in various disease areas. Visual inspection of these signatures readily highlights some patterns. For example, there is a specific group of side effects (E3) associated with anti-infective drugs, and ophthalmic drugs have similar mechanisms of action (B1) but varied chemistries (A1). We found that CC similarity searches greatly increase the chance of identifying drug properties compared to chemical similarities alone¹⁹ (Fig. 2c), partly because individual CC spaces have

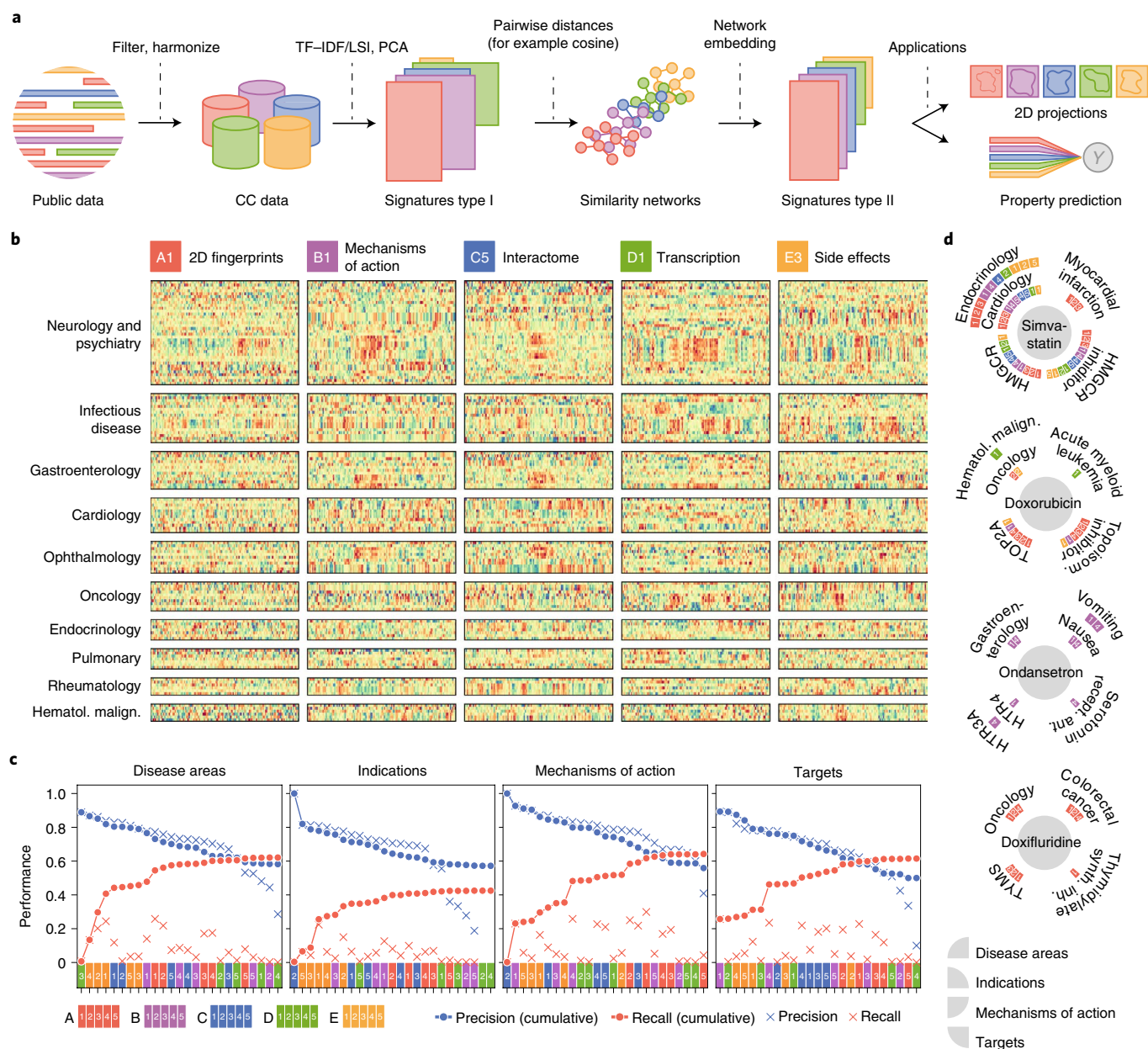
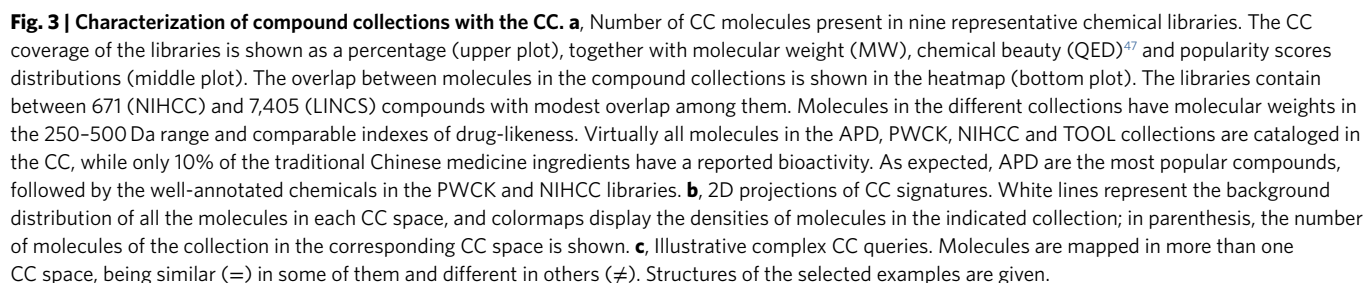


Fig. 2 | CC signatures visualized. **a**, Scheme of the CC pipeline. Public data are filtered, harmonized and unified in the 5×5 CC organization. For each CC space, we obtain type I signatures by doing a (TF-IDF) LSI/PCA dimensionality reduction. With signatures type I, molecules can be compared pairwise to obtain a similarity network. A network embedding algorithm (node2vec) is then applied to derive fixed-length signatures (type II). **Type I and/or type II signatures can be used for customary machine learning tasks such as data visualization and property prediction.** **b**, We plot the numerical values of type II signatures for drugs extracted from the Drug Repurposing Hub¹⁹, and organize them by disease areas. We chose one illustrative dataset for each CC level; namely, A1, B1, C5, D1 and E3. Signatures show, for instance, how chemically unrelated neurological drugs elicit similar patterns of side effects. Likewise, ophthalmological drugs sharing mechanism of action trigger different transcriptional responses. **c**, Precision and recall of label predictions (disease areas, indications, mechanisms of action and targets from the Drug Repurposing Hub, Methods). CC spaces are sorted by precision (blue). Recall of molecule-label pairs is shown in red. Dots correspond to cumulative performances (that is, appending molecule-label pairs predicted by CC spaces consecutively). Crosses denote individual performances of CC spaces. **d**, Examples of true positives, indicating the CC spaces that account for the prediction. Please note that the Drug Repurposing Hub was not included in the CC at the time of compilation. TF-IDF, term frequency-inverse document frequency transformation; hematol. malign., hematological malignancy; topoisom., topoisomerase; recept. ant., receptor antagonist; synth. inh., synthase inhibitor.

incomplete drug coverage (Supplementary Fig. 8a), and partly because different types of CC data capture different kinds of similarities between drugs (Supplementary Fig. 8b). For example, several CC spaces simultaneously accounted for the relationship between simvastatin, HMGCR inhibition and myocardial infarction (Fig. 2d). In the case of doxorubicin, its capacity to inhibit TOP2A was captured by chemical features, while its association with acute

myeloid leukemia was identified using transcriptional signatures (D1). For other drugs, such as ondansetron, the association with gastroenterology was more trivially provided by target annotations already present in the CC (B1 and B4), whereas for some drugs (for example, doxifluridine), the chemical similarity to other well-annotated compounds was enough to correctly uncover main therapeutic properties.



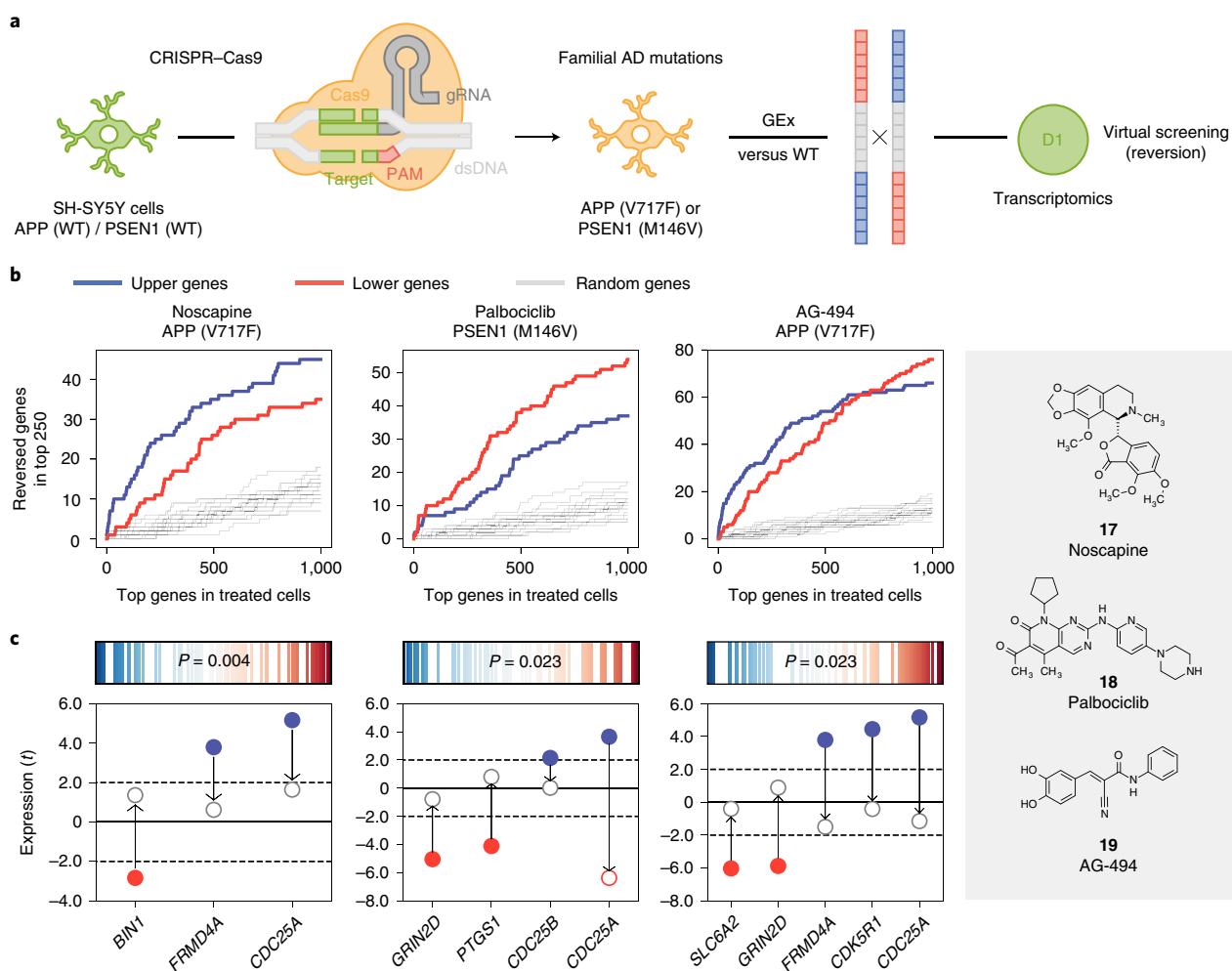


Fig. 4 | Signature reversion of Alzheimer's disease-specific transcriptional profiles. **a**, Scheme of the methodology. SH-SY5Y cells were modified with CRISPR to harbor fAD mutations. Alzheimer's disease (AD)-specific transcriptional signatures were obtained by differential gene expression analysis of mutated-versus-WT gene expression profiles. These signatures were flipped (reversed) and converted to the D1 CC format. Drug candidates were selected based on D1 similarities to the signatures. **b**, Experimental results for the three tested candidates, namely noscapine (**17**), palbociclib (**18**) and AG-494 (**19**). In the x axis, genes are ranked by differential gene expression of treated-versus-untreated mutated cells (APP^{V717F} or PSEN1^{M146V}); this axis relates to both tails of the ranked list (up/down). Correspondingly, in the y axis we count the number of genes in the mutated-versus-WT signatures that were reverted on treatment (top 250 genes, up- (blue) and down- (red) regulations). For example, ~20 of the up-regulated (blue) genes in PSEN1^{M146V} cells are in the top 500 down-regulated genes after treatment with palbociclib, and ~40 of the down-regulated (red) genes in the PSEN1^{M146V}-versus-WT comparison are among the top 500 up-regulated genes when these mutated cells are treated with palbociclib. **c**, Reversion of Alzheimer's disease-related genes. The upper plots show the tendency of Alzheimer's disease genes (according to OpenTargets) to have extreme reversion scores. Reversion scores measure the ratio between ranks in the mutated-versus-WT signatures and flipped (reversed) ranks on treatment of the mutated cells with the drug. Blue (left of the axis) denotes genes that were up-regulated in the mutated-versus-WT signature and down-regulated on treatment, and red (right of the axis) denotes genes that were down-regulated in mutated cells and up-regulated on treatment. The P value is calculated with a weighted one-sided Kolmogorov-Smirnov test based on the absolute value of these reversion scores, that is, it measures the 'extremity' of Alzheimer's disease genes. In the bottom plots, we focus on Alzheimer's disease genes that were up- (blue) and down- (red) regulated (t score) in the mutated-versus-WT comparison (bold dots), and we show their expression in the treated-versus-WT comparison (empty dots). Three independent experiments ($n=3$) were performed in all the experiments shown. GEx, gene expression.

Visualizing collections of compounds. CC signatures can be projected to two dimensions, providing new insights into compound libraries (Fig. 3a,b). For instance, Fig. 3b shows that, compared to preclinical libraries, APD map to a limited area of the physicochemical parameter space (A5), and reveals the structural diversity of screening libraries (A4). Experimental drugs are shown to address mechanisms of action (B1) not covered by APD, metabolites or tool compounds. Likewise, we see how they can elicit new transcriptional changes (D1) and how natural products, such as traditional Chinese medicines, may offer new possibilities. We can also observe that a diverse compound collection (Prestwick Chemical Library,

PWCK) may trigger a limited set of morphological changes (D4). In the clinical categories, we see the different zones painted by experimental drugs and traditional Chinese medicines (E2), and we also observe differences in the disease landscapes of endogenous and exogenous compounds (E4).

Further, combining 2D plots throughout the CC facilitates a better understanding of subgroups of compounds, and may inspire complex queries to identify molecules that fulfill multiple characteristics (Fig. 3c). For instance, despite being structurally diverse (A4), antitumor compounds chlorambucil (**5**), mitomycin C (**6**) and teniposide (**7**) trigger similar transcriptional responses (D1)

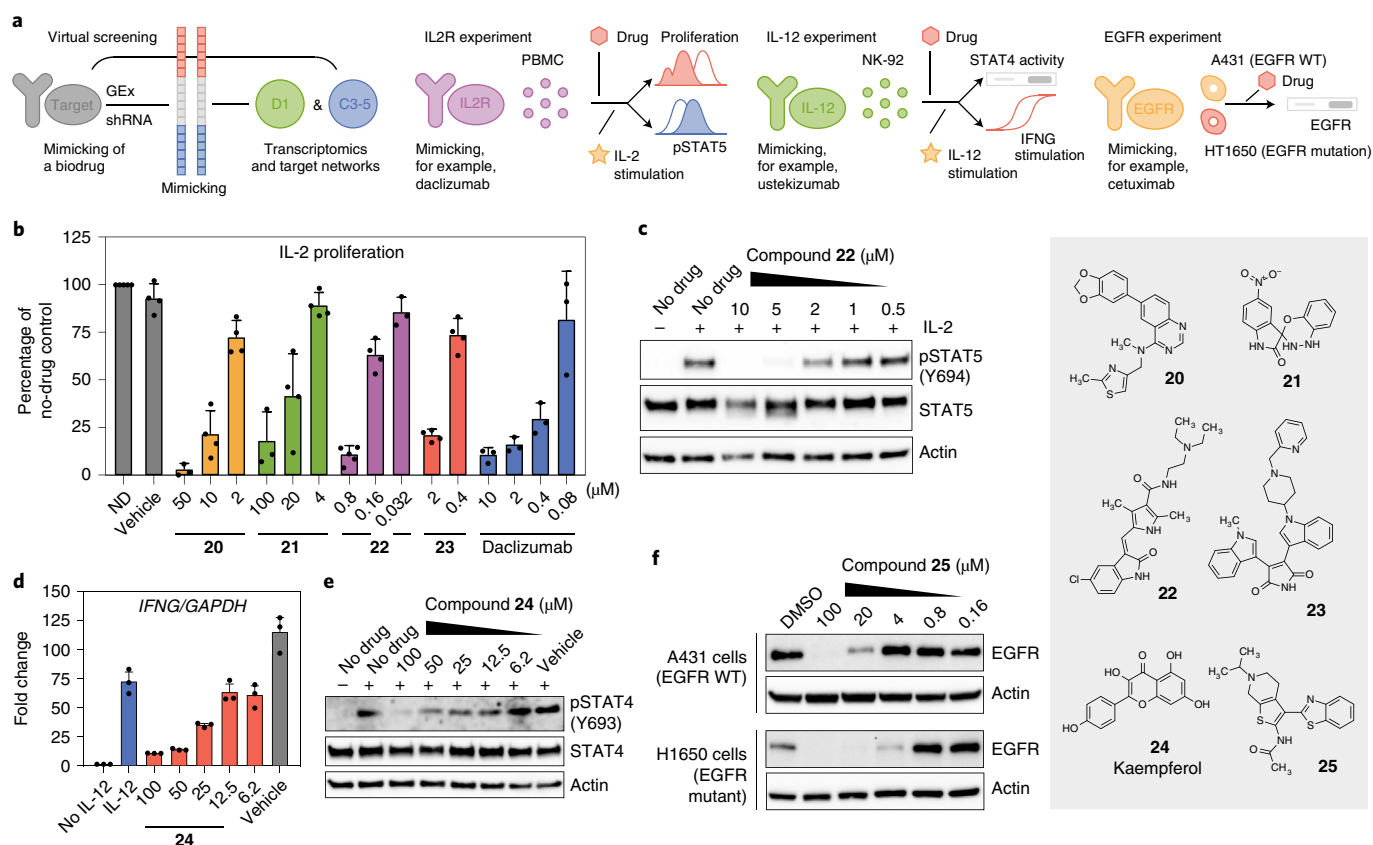


Fig. 5 | Discovery of chemical analogs of biologics. **a**, Scheme of the methodology. We look for compounds whose gene expression signatures (D1) would mimic gene expression signatures corresponding to the shRNA knock-down of the target of interest. In addition, we do a networks-level (C3–5) signature matching of the target profiles with those of the compounds. Candidates for IL-2 receptor, IL-12 and EGFR are tested in different experimental setups. **b**, CD3/CD28 prestimulated PBMCs were left without treatment for 3 d, labeled with CFSE and then stimulated with IL-2 (0.5 ng ml⁻¹) in the presence of the indicated compounds. Three days after stimulation, proliferation was measured by flow cytometry as CFSE label decay and normalized compared to the cells stimulated in the absence of drug (ND). Mean \pm s.d. of 3–5 independent experiments are shown, as illustrated by the dots in each barplot. **c**, IL-2-induced STAT5 phosphorylation in PBMCs quantified by western blot for compound 22. One representative experiment is shown ($n=3$). **d**, NK-92 cells were stimulated with IL-12 (50 ng ml⁻¹) in the presence of the indicated concentration of compound 24 (kaempferol). *IFNG* messenger RNA levels after 6 h were quantified by qPCR. Mean \pm s.d. of three independent experiments are shown. **e**, Phosphorylation of STAT4 at tyrosine 693 was assessed by western blot 1 h after stimulation with IL-12. Total STAT4 and actin antibodies were used as controls. One representative experiment is shown ($n=3$). **f**, A431 and H1650 cells were treated for 24 h with the indicated concentrations of compound 25 (APE1 inhibitor III). We quantified EGFR protein by western blot. Actin was used as a loading control. Representative blots out of three independent experiments are shown. GEx, gene expression.

and show similar cell-sensitivity profiles (D2), consistent with their known capacity to induce DNA damage—an uncharacterized compound (8) was found in this subgroup. We also identified a group of broad-spectrum CDK inhibitors (9, 10, 11 and 12) that induce a precise transcriptional response (D1). Conversely, we noticed that compounds within antibiotic classes (for example, beta-lactams 13 or sulfonamides 14) may be transcriptionally diverse in human cells (D1). Finally, we found compounds (15, 16) targeting kinases (B4) in various signaling pathways (mTOR/PIK3CA and Raf1/MAP2K1/MAP2K2, respectively) that are close in the interactome space (C5) and induce similar cell responses (D1), in agreement with a reported pathway cross-talk with potential for combination therapies²⁰.

Reversion of Alzheimer’s disease signatures. Having demonstrated the value of CC signatures to broadly characterize compound collections, we sought to explore their capacity to enable computational tasks that cannot be achieved using chemical information alone. A unique feature of CC signatures is that they can be matched to disease and genetic omics data. For instance, comparison of gene expression signatures in cells can reveal compounds

that ‘revert’ transcriptional disease signatures^{21,22}. Typically, these studies require intensive preprocessing²³, since direct comparisons of gene expression profiles, even within replicates, show modest correlations, and cell-specific biases can confound the analyses²². The CC pipeline handles the issues related to multiple doses, time points and cell lines, and returns only one D1 signature per compound (Supplementary Fig. 1).

The capacity of drugs to revert cancer gene expression profiles correlates with their efficacy²⁴ and, indeed, using D1 signatures we obtained similar results on the Genomics of Drug Sensitivity in Cancer (GDSC) panel of cancer cell lines²⁵ (Supplementary Fig. 9). The CC spaces are enriched in data obtained from tumor cell lines, and to evaluate the power of D1 signatures outside the realm of cancer, we engineered new cells for which no perturbation experiments are available. To this end, we developed cellular models of Alzheimer’s disease by introducing familial Alzheimer’s disease (fAD) mutations into human SH-SY5Y cells, which are known to recapitulate phenotypes related to neurodegenerative disorders²⁶. Using CRISPR–Cas9-induced homology-directed repair (HDR), we obtained clones harboring the fAD PSEN1^{M146V} or the APP^{V717F} mutations (Methods and Supplementary Fig. 10a,b). As expected, engineered cells showed

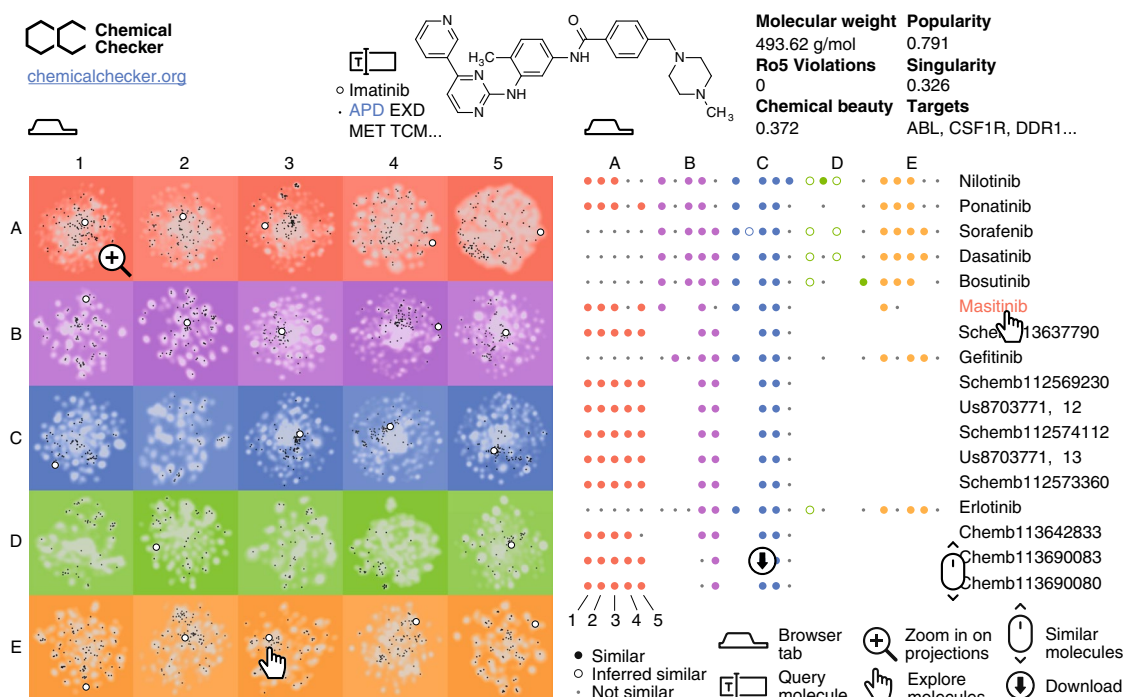


Fig. 6 | Representation of the CCweb resource. The left tab (home page) is an interactive panel of 2D projections, where the query molecule (for example, imatinib, white dot) can be compared to the CC background (in gray) and to other molecules of interest such as APD (in black). The right tab (exploration page) displays molecules that are similar to the query one. Similarities are measured across the 25 CC spaces (A1–E5).

an increased extracellular ratio of amyloid β (A β) 42 to A β 40 (Supplementary Fig. 10c), which is a hallmark of fAD mutations²⁷.

We measured the transcriptional signatures of PSEN1^{M146V}-versus-wild type (WT) and APP^{V717F}-versus-WT cells, which we flipped (that is, converting up- to down-regulated genes and vice versa) and adapted to the CC format (Fig. 4a and Methods). Then we simply ran a similarity search between the CC signatures of the compounds available in D1 and the reverted Alzheimer's disease-specific signatures (Supplementary Data 1). We identified 35 chemically diverse compounds that might have the potential to cancel out transcriptional traits of fAD mutations (Supplementary Fig. 11 and Supplementary Data 1). Of these, three, namely nescapine (17) (for the reversion of APP^{V717F} signature), palbociclib (18) (for the reversion of PSEN1^{M146V} signature) and the epidermal growth factor receptor (EGFR) inhibitor AG-494 (19) (for the reversion APP^{V717F} signature), showed an effect on the secretion of A β 40 and A β 42 in SH-SY5Y cells (Supplementary Fig. 10d).

We confirmed that genes up-regulated in SH-SY5Y fAD mutants were indeed down-regulated on treatment with the drugs, and vice versa (Fig. 4b). Moreover, the three drug treatments substantially reverted a subset of genes strongly linked to Alzheimer's disease²⁸ (Fig. 4c), including the recovery of the expression levels of *GRIN2D*, a glutamate receptor involved in synaptic transmission²⁹ and *BIN1*, a gene involved in synaptic vesicle endocytosis and strongly associated with Alzheimer's disease risk³⁰.

Mimicking the activity of biologics against IL2R, IL-12 and EGFR. Biologics are a family of medicines that includes antibodies and recombinant proteins. Although expensive and prone to pharmacokinetic issues³¹, biologics have the advantage that they bind with high specificity to their targets, which may be proteins considered undruggable by small molecules. We hypothesized that CC signatures could be used to find compounds that match the effect of biologics. Moreover, signatures corresponding to other spaces present in CC (for example, C3–C5) could be used to filter potential

hits. Thus, we devised a strategy that exploits the signature matching capacity of the CC and identifies compounds that could mimic the gene expression profile induced by certain biologics (D1), possibly via alternative targets participating in related biological processes (C3–C5). After an exploratory analysis (Supplementary Data 2 and Methods) we selected three biologic targets, namely the interleukin (IL)-2 receptor (IL2R), IL-12 and EGFR (Fig. 5a), based on the public availability of short-hairpin RNA interference (knock-down) experiments³².

Daclizumab is a monoclonal antibody targeting the alpha subunit of IL2R, and it is approved for the prevention of transplant rejection. Our computational search highlighted 23 diverse compounds that might mimic daclizumab (Supplementary Fig. 12 and Supplementary Data 2). We were able to purchase 19 of these compounds, and we tested their effect in the proliferation of primary human peripheral blood mononuclear cells (PBMC) stimulated with IL-2 (ref. 33). Fourteen significantly inhibited PBMC proliferation without substantial effects on cell viability (Supplementary Fig. 13); 13 of the 14 also significantly inhibited PHA-stimulated proliferation³² (Supplementary Table 2 and Supplementary Fig. 14). The hit rate of comparable high-throughput assays is 0.5–15% (PubChem BioAssays AIDs: 371, 463, 575, 598, 648, 719, 772 and 2303). In (partially) IL-2 independent cells, the antiproliferative effect was only moderate (Supplementary Fig. 15). Figure 5b shows confirmatory dose–response curves for four of the candidates, including previously uncharacterized compounds (20 and 21). Further analysis revealed that compound 22 inhibited STAT5 phosphorylation on IL-2 stimulation (Fig. 5c), indicating it acts in the same signaling pathway as daclizumab. On the contrary, compounds 20, 21 and 23 did not block STAT5 phosphorylation, suggesting that their antiproliferative effect blocks a complementary pathway (Supplementary Fig. 16).

Ustekinumab, a monoclonal antibody targeting IL-12 and IL-23 interleukins, is approved for the treatment of psoriasis and has potential in autoimmune syndromes³³. It blocks interferon-gamma

(IFNG) production from natural killer (NK) cells that is induced by IL-12 and IL-23 receptor binding and STAT4 phosphorylation³⁴. Our search for compounds that can match C3–5 and D1 signatures of ustekinumab highlighted 17 candidates (Supplementary Data 2 and Supplementary Fig. 12). We tested the capacity of 11 of them to block IL-12-induced IFNG production in natural killer (NK) cells. One of the compounds, kaempferol (**24**), inhibited *IFNG* transcription in a dose-dependent manner (Fig. 5d). Moreover, kaempferol inhibited the phosphorylation of STAT4 at tyrosine 693 in response to IL-12, indicating that this compound exerts its action in an early step of IL-12 signaling (Fig. 5e).

Monoclonal antibodies targeting EGFR (for example, cetuximab) are used to treat colon and head and neck cancers³⁵. Our CC signature matching search highlighted three candidates (Supplementary Data 2), including apigenin and tanespmycin (17-AAG), which are known to affect EGFR signaling in vitro and in vivo, and to synergize with cetuximab^{36–38}. The third compound was an apurinic/aprimidinic endodeoxyribonuclease (APE1) inhibitor (**25**) that, to our knowledge, has no reported connection to EGFR. Treatment with compound **25** degraded EGFR in a dose-dependent manner in WT and EGFR^{ΔE746–A750} mutated cells (Fig. 5f).

Similarity searches in the CC. We built a web-based resource (CCweb, at <https://chemicalchecker.org>) to facilitate access to our data. As shown in Fig. 6, the CCweb displays the 2D projection of each dataset, offering the possibility to use chemical libraries as landmark points and highlighting how individual compounds are distributed and related to each other. In addition, it provides ‘popularity’ and ‘singularity’ (Fig. 1d) scores for all compounds, which account for the number of CC spaces related to a certain molecule and the uniqueness (dissimilarity) of a molecule with respect to the rest of compounds, respectively. Moreover, given a molecule of interest, the CCweb retrieves similar molecules in all 25 CC spaces. Most small-molecule search engines available to the community are based on chemical similarities, whereas CCweb offers search capacity based on biological similarities. CC signatures can be downloaded from the CCweb or simply accessed via a representational state transfer application programming interface. The entire CCweb resource, including the underlying data and signatures, will be updated every 6 months. There is a link to the full code of our resource in the CCweb page.

Discussion

As small-molecule bioactivity data continue to grow in size and diversity, it is essential to present them in a format accessible to most researchers. Big initiatives such as OpenPHACTS³⁹ and Illuminating the Druggable Genome (IDG)⁴⁰ are undertaking this task, storing links between compounds, genes and diseases in a relational scheme that is ideal for browsing and formulating mechanistic hypotheses. With the CC, we propose an alternative framework based on chemical and biological signatures of compounds. CC signatures are numeric vectors that embed information of a given type (for example, binding experiments, cell-sensitivity profiles or drug side effects) and are suitable for similarity measurements, clustering, visualization and prediction tasks. Such capabilities, we believe, are essential to bridge the gap between relational databases and frontline machine learning algorithms that are able to handle millions of samples but require input data to be expressed in vector format.

The signature-based representation of compounds pushes the similarity principle beyond chemical properties to various ambits of biology. For instance, our preliminary experiments identified candidates to revert Alzheimer’s disease transcriptional signatures, and we devised a strategy to propose small-molecule mimetics of biologics. We also used signatures based on pathways, biological processes and networks to gain confidence in our predictions. More generally, we have visualized compound collections by mapping

them to different bioactivity spaces, and have shown that similarity searches inside the CC recapitulate drug indications and mechanisms of action.

Other applications of the CC include the replacement of traditional chemical fingerprints with CC signatures in supervised machine learning tasks such as ligand-based target prediction⁴¹, as well as large-scale unsupervised predictions against disease profiles^{18,42} based on the notion of signature connectivity. Further, recent advances in machine learning suggest that a signature-guided de novo design of small molecules is possible⁴³, offering an opportunity to further populate the bioactive chemical space.

The current version of the CC contains ~800,000 molecules. All of them have experimental annotations in at least one of the biological levels (B–E, typically B4 and B5). The known chemical space is much larger than this, containing millions of commercial compounds and a cosmic number of synthetically accessible virtual molecules⁴⁴. A good proportion of the molecules will not be bioactive, falling outside the scope of the CC. However, the bioactive chemical space remains mostly uncharted⁴⁵, meaning that the current CC data are incomplete, especially for the higher-order (phenotypic and clinical) layers. We have observed remarkable correlations between the different data types contained in the CC, which suggests that inference of CC signatures would be possible for poorly characterized compounds. Future directions for the CC include the massive prediction of missing bioactivity data based on the currently assembled resource, offering a means to rapidly characterize any molecule of interest. Likewise, we expect the CC to evolve in terms of data types as new screening technologies continue to emerge. As the CC grows in complexity, large-scale data fusion algorithms⁴⁶ will be instrumental to enable a global view of the similarity space and ensure that the simple, convenient organization of the resource is maintained.

Online content

Any methods, additional references, Nature Research reporting summaries, source data, extended data, supplementary information, acknowledgements, peer review information; details of author contributions and competing interests; and statements of data and code availability are available at <https://doi.org/10.1038/s41587-020-0502-7>.

Received: 23 August 2019; Accepted: 27 March 2020;

Published online: 18 May 2020

References

- Sterling, T. & Irwin, J. J. ZINC 15—ligand discovery for everyone. *J. Chem. Inform. Model.* **55**, 2324–2337 (2015).
- Gaulton, A. et al. The ChEMBL database in 2017. *Nucleic Acids Res.* **45**, D945–D954 (2017).
- Wang, Y. et al. PubChem BioAssay: 2017 update. *Nucleic Acids Res.* **45**, D955–D963 (2017).
- Wishart, D. S. Chapter 3: small molecules and disease. *PLOS Comput. Biol.* **8**, e1002805 (2012).
- Duran-Frigola, M., Rossell, D. & Aloy, P. A chemo-centric view of human health and disease. *Nature Commun.* **5**, 5676 (2014).
- Rouillard, A. D. et al. The harmonizome: a collection of processed datasets gathered to serve and mine knowledge about genes and proteins. *Database* **2016**, baw100–baw100 (2016).
- Newman, D. J. & Cragg, G. M. Natural products as sources of new drugs from 1981 to 2014. *J. Nat. Prod.* **79**, 629–661 (2016).
- Rodrigues, T., Reker, D., Schneider, P. & Schneider, G. Counting on natural products for drug design. *Nat. Chem.* **8**, 531–541 (2016).
- Welsch, M. E., Snyder, S. A. & Stockwell, B. R. Privileged scaffolds for library design and drug discovery. *Curr. Opin. Chem. Biol.* **14**, 347–361 (2010).
- Bleicher, K. H., Böhm, H.-J., Müller, K. & Alanine, A. I. Hit and lead generation: beyond high-throughput screening. *Nat. Rev. Drug Disc.* **2**, 369–378 (2003).
- Holbeck, S. L., Collins, J. M. & Doroshow, J. H. Analysis of food and drug administration-approved anticancer agents in the NCI60 panel of human tumor cell lines. *Mol. Cancer Therap.* **9**, 1451–1460 (2010).

12. Seashore-Ludlow, B. et al. Harnessing connectivity in a large-scale small-molecule sensitivity dataset. *Cancer Discov.* **5**, 1210–1223 (2015).
13. Campillos, M., Kuhn, M., Gavin, A.-C., Jensen, L. J. & Bork, P. Drug target identification using side-effect similarity. *Science* **321**, 263–266 (2008).
14. Petrone, P. M. et al. Rethinking molecular similarity: comparing compounds on the basis of biological activity. *ACS Chem. Biol.* **7**, 1399–1409 (2012).
15. Papadatos, G., Gaulton, A., Hersey, A. & Overington, J. P. Activity, assay and target data curation and quality in the ChEMBL database. *J. Comput. Aided Mol. Des.* **29**, 885–896 (2015).
16. Duran-Frigola, M., Mateo, L. & Aloy, P. Drug repositioning beyond the low-hanging fruits. *Curr. Opin. Syst. Biol.* **3**, 95–102 (2017).
17. Nguyen, D. T. et al. Pharos: collating protein information to shed light on the druggable genome. *Nucleic Acids Res.* **45**, D995–D1002 (2017).
18. Duran-Frigola, M., Fernandez-Torras, A., Bertoni, M. & Aloy, P. Formatting biological big data for modern machine learning in drug discovery. *WIREs Comp. Mol. Sci.* **9**, e1408 (2018).
19. Corsello, S. M. et al. The Drug Repurposing Hub: a next-generation drug library and information resource. *Nat. Med.* **23**, 405–408 (2017).
20. Jokinen, E. & Koivunen, J. P. MEK and PI3K inhibition in solid tumors: rationale and evidence to date. *Ther. Adv. Med. Oncol.* **7**, 170–180 (2015).
21. Lamb, J. et al. The connectivity map: using gene-expression signatures to connect small molecules, genes, and disease. *Science* **313**, 1929–1935 (2006).
22. Subramanian, A. et al. A next generation connectivity map: L1000 platform and the first 1,000,000 profiles. *Cell* **171**, 1437–1452 (2017).
23. Filzen, T. M., Kutchukian, P. S., Hermes, J. D., Li, J. & Tudor, M. Representing high throughput expression profiles via perturbation barcodes reveals compound targets. *PLoS Comput. Biol.* **13**, e1005335 (2017).
24. Chen, B. et al. Reversal of cancer gene expression correlates with drug efficacy and reveals therapeutic targets. *Nat. Commun.* **8**, 16022 (2017).
25. Iorio, F. et al. A landscape of pharmacogenomic interactions in cancer. *Cell* **166**, 740–754 (2016).
26. Encinas, M. et al. Sequential treatment of SH-SY5Y cells with retinoic acid and brain-derived neurotrophic factor gives rise to fully differentiated, neurotrophic factor-dependent, human neuron-like cells. *J. Neurochem.* **75**, 991–1003 (2000).
27. Tanzi, R. E. The genetics of Alzheimer disease. *Cold Spring Harb. Perspect. Med.* **2**, a006296 (2012).
28. Carvalho-Silva, D. et al. Open Targets Platform: new developments and updates two years on. *Nucleic Acids Res.* **47**, D1056–D1065 (2019).
29. Perszyk, R. E. et al. GluN2D-containing N-methyl-D-aspartate receptors mediate synaptic transmission in hippocampal interneurons and regulate interneuron activity. *Mol. Pharmacol.* **90**, 689–702 (2016).
30. Harold, D. et al. Genome-wide association study identifies variants at CLU and PICALM associated with Alzheimer's disease. *Nat. Genet.* **41**, 1088–1093 (2009).
31. Anselmo, A. C., Gokarn, Y. & Mitragotri, S. Non-invasive delivery strategies for biologics. *Nat. Rev. Drug Discov.* **18**, 19–40 (2018).
32. Depper, J. M., Leonard, W. J., Robb, R. J., Waldmann, T. A. & Greene, W. C. Blockade of the interleukin-2 receptor by anti-Tac antibody: inhibition of human lymphocyte activation. *J. Immunol.* **131**, 690–696 (1983).
33. Benson, J. M. et al. Therapeutic targeting of the IL-12/23 pathways: generation and characterization of ustekinumab. *Nat. Biotechnol.* **29**, 615–624 (2011).
34. Reddy, M. et al. Modulation of CLA, IL-12R, CD40L, and IL-2R α expression and inhibition of IL-12- and IL-23-induced cytokine secretion by CNO 1275. *Cell Immunol.* **247**, 1–11 (2007).
35. Xu, M. J., Johnson, D. E. & Grandis, J. R. EGFR-targeted therapies in the post-genomic era. *Cancer Metastasis Rev.* **36**, 463–473 (2017).
36. Masuelli, L. et al. Apigenin induces apoptosis and impairs head and neck carcinomas EGFR/ErbB2 signaling. *Front. Biosci.* **16**, 1060–1068 (2011).
37. Hu, W. J., Liu, J., Zhong, L. K. & Wang, J. Apigenin enhances the antitumor effects of cetuximab in nasopharyngeal carcinoma by inhibiting EGFR signaling. *Biomed. Pharmacother.* **102**, 681–688 (2018).
38. Sawai, A. et al. Inhibition of Hsp90 down-regulates mutant epidermal growth factor receptor (EGFR) expression and sensitizes EGFR mutant tumors to paclitaxel. *Cancer Res.* **68**, 589–596 (2008).
39. Williams, A. J. et al. Open PHACTS: semantic interoperability for drug discovery. *Drug Disc. Today* **17**, 1188–1198 (2012).
40. Rodgers, G. et al. Glimmers in illuminating the druggable genome. *Nat. Rev. Drug Disc.* **17**, 301–302 (2018).
41. Wu, Z. et al. MoleculeNet: a benchmark for molecular machine learning. *Chem. Sci.* **9**, 513–530 (2018).
42. Lee, Y. S. et al. A computational framework for genome-wide characterization of the human disease landscape. *Cell Syst.* **8**, 152–162 (2019).
43. Mendez-Lucio, O., Baillif, B., Clevert, D. A., Rouquie, D. & Wichard, J. De novo generation of hit-like molecules from gene expression signatures using artificial intelligence. *Nat. Commun.* **11**, 10 (2020).
44. Raymond, J.-L. The Chemical Space Project. *Acc. Chem. Res.* **48**, 722–730 (2015).
45. Irwin, J. J., Gaskins, G., Sterling, T., Mysinger, M. M. & Keiser, M. J. Predicted biological activity of purchasable chemical space. *J. Chem. Info. Modeling* **58**, 148–164 (2018).
46. Wang, B. et al. Similarity network fusion for aggregating data types on a genomic scale. *Nat. Methods* **11**, 333–337 (2014).
47. Bickerton, G. R., Paolini, G. V., Besnard, J., Muresan, S. & Hopkins, A. L. Quantifying the chemical beauty of drugs. *Nat. Chem.* **4**, 90–98 (2012).

Publisher's note Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.

© The Author(s), under exclusive licence to Springer Nature America, Inc. 2020

Methods

Raw data. Small-molecule entries were collected from several resources (Table 1 and Supplementary Fig. 1) and stored by standard InChIKey. The InChIKey is a 25-character string that encodes the connectivity of the molecule (first 14 characters), other details such as stereochemistry (next eight characters), the kind and version of the key (next two characters), and the protonation state (last character). To assign an InChIKey to each small molecule, we read the structure as given by the source database (usually a SMILES string) and followed a standardization procedure consisting of salt and solvent removal, charge neutralization and the application of rules to tautomeric groups (<https://github.com/flatiron/standardiser>).

A: chemistry. A1: 2D fingerprints. The 2,048-bit Morgan fingerprints (radius of 2) were calculated using the RDKit (<http://rdkit.org>).

A2: 3D fingerprints. The 1,024-bit E3FP fingerprints (<https://github.com/keiserlab/e3fp>) were calculated by merging the results of the three best conformers obtained with a UFF energy minimization, as recommended in the E3FP publication⁴⁸.

A3: scaffolds. We extracted the Murcko's scaffold of each molecule⁴⁹. In addition, we derived the molecular framework of the scaffold, that is, all heavy atoms were converted to carbon atoms and all bonds were simplified to single bonds. When no scaffold could be obtained, we kept the full structure of the molecule and the corresponding framework. The 1,024-bit Morgan fingerprints (radius of 2) were then calculated for each molecule and concatenated in a 2,048-bit fingerprint.

A4: structural keys. The widely used, human-readable molecular access system 166 keys⁵⁰ were calculated using the RDKit. Molecular access system keys represent structural features relevant to medicinal chemistry. Each key is associated to a SMARTS pattern. Although more fine-grained fingerprints (for example, A1) are in general preferred in modern chemoinformatics tasks, we found the coarser A4 fingerprints to be convenient for global exploration task such as 2D projections and visualization (Fig. 3b).

A5: physicochemical parameters. For each molecule, we calculated the molecular weight, number of heavy atoms, number of heteroatoms, number of rings, number of aliphatic rings, number of aromatic rings, number of hydrogen bond acceptors, number of hydrogen bond donors and number of rotatable bonds. We predicted logP, molecular refractivity, and polar surface area using RDKit. In addition, we flagged the structural alerts proposed by Hopkins and coworkers⁴⁷ and those listed in ChEMBL² (v.22, <https://www.ebi.ac.uk/chembl>). We also counted Lipinski's rule-of-five violations⁵¹ and rule-of-three violations⁵². Finally, the chemical beauty (QED)⁴⁷ was quantified using the Silicos-IT kit (<http://silicos-it.be.s3-website-eu-west-1.amazonaws.com/>).

B: targets. B1: mechanism of action. Mechanisms of action of approved and experimental drugs were collected from DrugBank⁵³ (v.4, <https://www.drugbank.ca>) by selecting those protein targets with a known pharmacological action and action mode. Similarly, we fetched from ChEMBL those drugs with a known mode of action. We distinguished between 'activation' modes (agonist, activator and so on) and 'inhibition' modes (antagonist, competitor and so on). Together with the identity of the protein targets, we retained protein class memberships (G protein-coupled receptors, kinases and so on) at all levels of the ChEMBL target hierarchy.

B2: metabolic genes. We collected drug-metabolizing enzymes, transporters and carriers from DrugBank. To these, we added proteins involved in drug metabolism as recorded in ChEMBL. As in B1, we retained protein class information.

B3: crystals. We downloaded ligand data from the Protein Data Bank (<https://www.rcsb.org>, February 2017). Protein structures bound to each small molecule were then annotated with family (F and T groups) and superfamily (H and X groups) information, following the Evolutionary Classification of Protein Domains⁵⁴ (ECOD v.1.4, <http://prodata.swmed.edu/ecod>).

B4: binding. Protein binding data were obtained from ChEMBL by searching for bioassays of 'binding' type, related to 'single proteins' with an experimental measure of standard type ('pChEMBL' value available). We also collected BindingDB records with activity expressed as concentrations⁵⁵ (<https://www.bindingdb.org>, February 2017). Data were discretized by applying the following activity cutoffs, recommended in Pharos⁵⁷ (<http://pharos.nih.gov/>): kinases ≤ 30 nM, G protein-coupled receptors ≤ 100 nM, nuclear receptors ≤ 100 nM, ion channels ≤ 10 μ M and others ≤ 1 μ M. We also kept activities one order of magnitude lower than the class-specific cutoff (to a maximum of 10 μ M), and gave these annotations half the weight in downstream analyses (that is, log₁₀ scaling). Finally, protein class hierarchy information was kept as in B1.

B5: high-throughput screening bioassays. The largest public repository of small-molecule screening data is PubChem Bioassays³. Bioactivity values from this repository were directly downloaded from ChEMBL, since the latter conveniently applies a processing pipeline that collects only confirmatory assays and maps related protein targets to UniProt identifiers. Most of the assays belong to the 'functional' category. For completeness, we included other functional assays available in ChEMBL. We chose a relaxed activity cutoff of 10 μ M, or checked for the word 'active' in the description of the assay. We kept the protein class hierarchy as in B1.

C: networks. C1: small-molecule roles. We downloaded the Chemical Entities of Biological Interest (ChEBI) ontology⁵⁶ (v.150, <http://www.ebi.ac.uk/chebi>). Only 'three-star' molecules were considered. The 'role' ontology was loaded as a directed graph capturing 'is a', 'is conjugate acid/base of', 'is enantiomer of', 'is tautomer of' and 'has role' relationships. In the ChEBI graph, molecules are 'leaves'. We searched for paths to reach the 'root' of the graph (that is, the 'role' node) from each of the leaves. Terms belonging to these paths were annotated to the corresponding molecules.

C2: metabolic pathways. We downloaded the reconstruction of human metabolism (Recon)⁵⁷ from Pathway Commons⁵⁸ (<http://www.pathwaycommons.org>, July 2017) in binary interaction form. Data were represented as an undirected graph where nodes are metabolites and edges denote reactions. We then computed an 'influence matrix' based on this metabolic network. In brief, positions in the influence matrix quantify the proximity between pairs of metabolites. The neighbors of each molecule were selected following a weighting scheme to favor proximal metabolites. Please see C5 for more details on how influence matrices are calculated and neighbors extracted and weighted therefrom.

C3: signaling pathways. The C2 space above is focused on endogenous metabolites. Conversely, this C3 space (and the C4 and C5 spaces) is aimed at any molecule with known protein targets. In this case, we list the biological pathways that may be affected by the interaction of a molecule with its targets. Human pathways were collected from Reactome⁵⁹ (<https://reactome.org>, May 2017), and we chose to use binding activities from B4, since this is an extensive dataset containing mostly literature data with well-accepted activity thresholds. In B4, 24.5% of the compound-protein interactions do not correspond to human proteins. These were mapped to their human orthologs using MetaPhOrs⁶⁰ (<http://orthology.phylomedb.org>, May 2017), following the observation that binding activities can be safely transferred between orthologous proteins⁶¹, especially if they belong to closely related species, as it is the case for B4 data⁶². Of all the nonhuman proteins mapped to the human orthologs, 94.4% were mammal proteins.

Molecules were annotated with Reactome pathways using a simple guilt-by-association approach; that is, a pathway was kept when at least one of its proteins was a target of the molecule. Pathways at all levels of the Reactome hierarchy were evaluated, and weight was given to the pathway annotation on the basis of the compound-target binding record (see B4).

C4: biological processes. We downloaded the Gene Ontology Annotation database (<https://www.ebi.ac.uk/GOA>, May 2017) and read the 'biological process' branch of the ontology as a directed acyclic graph ('is a' relationships). Proteins were annotated with their Gene Ontology Annotation biological process terms plus parent terms (up to the root of the directed acyclic graph). Similar to C3, we associated molecules with biological process terms by simply checking the annotations of the molecule targets (B4).

C5: interactomes. We collected five representative PPI networks, namely STRING (score of >700, that is, high confidence)⁶³ (v.10, <https://string-db.org>) (14,725 proteins (*p*), 300,686 interactions (*i*)), InWeb (score of 0.5)⁶⁴ (<http://www.intomics.com/inbio/map>, March 2017) (10,100 *p*, 168,970 *i*), a portion of Pathway Commons containing interactions from known pathways (Kyoto Encyclopedia of Genes and Genomes⁶⁵, NetPath⁶⁶, PANTHER⁶⁷ and WikiPathways⁶⁸) (9,344 *p*, 242,962 *i*), an in-house network of physical binary PPIs⁶⁹ (13,038 *p*, 64,659 *i*), and a network of metabolic genes based on Recon (v.2, <http://vmh.uni.lu>) (1,628 *p*, 246,937 *i*). To build this last network, we linked two metabolic proteins (enzymes or transporters) when the product metabolite of the first was the substrate of the second, or when both were needed to perform a certain reaction, suggesting that they are part of the same protein complex. Edges between proteins were weighted inversely proportional to the number of reactions involving their shared metabolites, so that 'currency' metabolites such as ATP and water had marginal impact on the network connectivity. To control for indirect associations, we deconvoluted the network⁷⁷ using edge weights and setting a network deconvolution score cutoff of 0.9.

The five networks above were treated separately in the following procedures. Given a PPI network containing *n* nodes, we computed a *n* × *n* 'influence matrix' using HotNet⁷⁰ (v.2, <https://github.com/raphael-group/hotnet2>). The influence matrix measures how likely a random walker departing from node *i* is to reach node *j*. This measure accounts for topological features such as centrality and betweenness, hence it is a more robust quantification of the relationship between nodes than the simple presence or absence of an interaction between them. Then, given the targets of a small molecule (B4), we looked in the matrix for the nodes that are most 'influenced' by these targets; that is, we retrieved proteins other than the target that are likely to be affected by the compound. The search for 'influenced' nodes was done as follows. First, nondiagonal values in the influence matrix were scaled from zero to ten and expressed as integers; as expected, most of the values were equal to 0, meaning that most proteins pairs were not influencing each other. Then, for each target of a certain compound, we kept proteins with a non-0 influence score (the target itself was given a score of ten and, when one protein was influenced by more than one target, the maximum score was kept). Finally, these scores were multiplied by the weight of the compound-target annotation (see B4). As a result, for each small molecule in each network, we obtained a weighted set of proteins that may be affected by the interplay with the targets. Results from the five different networks were concatenated for further analyses.

D: cells. D1: gene expression. Transcriptional profiles of treated cultured cells were obtained from the L1000 Connectivity Map²² (Phase I, GSE92742 and Phase II, GSE70138 in the Gene Expression Omnibus, March 2017). In this dataset, each ‘perturbagen’ (small molecule, shRNA or overexpressed gene) has several gene expression signatures assigned, corresponding to different doses, times of exposure, cell lines and so on. We took level 5 (replica-aggregated) signatures, considering both landmark and inferred gene expressions. Signatures with a low correlation between replicates ($\text{‘distil_cc_q75’} < 0.2$) were discarded. Following the authors’ recommendations²², we picked an ‘exemplar’ signature for each perturbagen in each available cell line by prioritizing signatures with a number of samples between two and six, and selecting the one with a highest transcriptional activity score. As a result, each perturbagen-cell-line pair has one (and only one) signature assigned.

After the filtering above, the complete L1000 Connectivity Map contained 22,118 perturbagens, each of them tested, on average, in 3.8 of 86 cell lines. A smaller, functionally diverse, and well-annotated subset of the data is the Touchstone dataset, which is focused on 8,880 perturbagens screened against a core collection of nine cell lines. The Touchstone dataset is the one that is queried in the online application of the L1000 Connectivity Map (<https://clue.io/l1000-query>), and we chose to use it as a reference collection of signatures. Accordingly, we measured pairwise similarities (connectivities) between the small-molecule signatures (trt_cp) of the full dataset (F) and the Touchstone (T) signatures (‘trt_cp’ , ‘trt_sh_cgs’ and ‘trt_oe’). To this end, we took the top 250 over- and underexpressed genes of the F signature⁷¹ and ran a two-way gene-set enrichment analysis (GSEA) against T signatures to obtain connectivity scores²² corresponding to the average between the GSEA enrichment score of up-regulated genes and the GSEA of down-regulated ones.

We normalized the connectivity scores (NCS) so that they were comparable between T-signature cell lines and perturbation types (small molecule, shRNA or gene overexpression). Normalization was simply done by dividing connectivity score by its average in each perturbation type category.

The CC is compound-centric, hence we summarized the results above, obtained for individual cell types, into a single measure of connectivity between F molecules and T perturbagens. A cell-summarized (consensus) connectivity score (NCS_{con}) was given by the maximum tertile statistic, first across T signatures and then across F signatures.

As a result, we obtained a F-versus-T connectivity matrix comparing the expression patterns of all molecules to the expression patterns of reference (Touchstone) perturbagens. Finally, we discretized this connectivity matrix by selecting, for each F molecule, significantly similar T perturbagens ($P < 0.01$, that is, 99% percentile of NCS_{con}). Molecules with fewer than five significantly similar T perturbagens were discarded.

D2: cancer cell lines. Modern cancer cell line panels such as the Cancer Cell Line Encyclopedia⁷², the GDSC²⁵ and the Cancer Therapeutics Response Portal⁷³ contain about a thousand cell lines but are short on screened molecules, having at most a few hundred of them. Conversely, the more classical NCI-60 cancer panel⁷⁴, while substantially narrower (60 cell lines), has almost 20,000 molecules screened for sensitivity, thus making it a better case for the CC. Indeed, Supplementary Fig. 17 shows that a relatively small number of cell lines is sufficient to accurately perform similarity searches across the D2 space. We collected z-transformed GI50 data from the National Institutes of Health (NIH) Developmental Therapeutics Program (<https://dtp.cancer.gov>, June 2016). Only molecules screened against at least 50 of the cell lines were considered. When more than one sensitivity profile was available for a given InChIKey, we kept the one with the largest number of assayed cell lines. This left us with a small-molecule sensitivity matrix that was 95.2% complete. Missing values were imputed using the MICE imputation algorithm over 100 iterations⁷⁵.

D3: chemical genetics. We downloaded chemical genetics data from MOSAIC⁷⁶ (<http://mosaic.cs.umn.edu>, September 2017). The raw chemical genetics dataset contains ~10,000 small molecules screened against ~300 yeast mutants. These ~300 yeast mutants were selected by the authors of the dataset so that they are representative of a broader panel of ~5,000 mutants. The ~10,000 × 300 chemical genetics matrix then becomes truly informative when it is compared to the ~5,000 × 300 genetic interaction matrix, in such a way that similarities between compounds and gene alterations can be discovered. This comparison is conveniently published in MOSAIC as a ‘gene target prediction’ file. We discretized the information in this file by keeping the identity of yeast mutants whose profiles had a similarity score above 7.12 (corresponding to a P value of 0.001) and, with half the weight, yeast mutants with a score above 3.37 ($P = 0.01$). Only 3,560 molecules passed this significance filtering.

D4: morphology. We downloaded the LDS-1195 dataset from the Library of Integrated Network-Based Cellular Signatures (LINCS) Data Portal (<http://lincsportal.ccs.miami.edu>), corresponding to cell painting morphological profiles⁷⁷. This dataset reports 812 cell image features measured after treatment of cells with ~30,000 compounds. To filter out molecules that do not have a substantial impact on cell morphology, we first counted the number of features (N_f) of each molecule that were significantly extreme ($P < 0.01$, that is, bottom 1% and top 99% of feature value distribution). We then repeated the same procedure to column-wise permuted versions of the data, and kept the N_{f0} point of $P < 0.01$ significance of this null distribution. Accordingly, we considered that molecules with $N_f < N_{f0}$ did not

trigger a significant morphological pattern, and we consequently discarded them; 12,075 molecules remained after the filtering.

D5: cell bioassays. We downloaded literature cell bioassay data from ChEMBL. We kept only standardized activity data given in commonly used units such as GI50, LC50 or IC50. Activities below 1 μM were retained, together with values beyond the 50% when data were percentual. We excluded cell lines that could not be mapped to the Cellosaurus ontology (v.22, <https://web.expasy.org/cellosaurus>). The Cellosaurus was used to identify and retain ‘derived from’ relationships between cell lines.

E: clinics. E1: therapeutic areas. We collected ATC classification system codes from DrugBank and the Kyoto Encyclopedia of Genes and Genomes. To capture the ATC hierarchy, we annotated molecules with their full ATC code (level 5), plus all higher levels (4 to 1).

E2: indications. We fetched approved and phase I-IV drug indications from ChEMBL and ReproDB⁷⁸ (v.1, <http://apps.chiragipgroup.org/repoDB>). ReproDB is an indication-oriented version of DrugBank. Unified Medical Language System disease terms in ReproDB were mapped to the MeSH vocabulary using DisGeNET⁷⁹ (v.4, <http://disgenet.org>) (MeSH is the preferred vocabulary in ChEMBL). We considered APD indications, together with those in clinical trials, from both databases. When a drug was indicated for more than one disease, weight was assigned to each indication depending on the clinical status (phase I to phase IV/approved), so that, for example, phase II annotations were twice as weighted as those of phase I. MeSH terms were spanned across the MeSH hierarchy as explained in E4. We kept the maximum weight for each parent term.

E3: side effects. We collected drug side effects from SIDER⁸⁰ (v.4, <http://sideeffects.embl.de>), expressed as Unified Medical Language System terms. We did not consider frequency information since we and others have found it to be too scarce for comprehensive statistical analyses^{81,82}.

E4: disease phenotypes. Associations between chemicals and disease phenotypes were downloaded from the Comparative Toxicogenomics Database (CTD)⁸³ (<http://ctdbase.org>, July 2016). We took only ‘curated’ CTD data. In CTD, compound-disease associations are classified as ‘therapeutic’ (T) or ‘marker/mechanism’ (M) (usually corresponding to a disease-causing effect). T and M annotations were kept separately for each molecule. CTD contains a medical vocabulary (MEDIC) that is essentially based on the MeSH hierarchy. For each annotated disease, we added parent terms all the way to the root of the MEDIC hierarchy.

E5: DDIs. To the best of our knowledge, DrugBank is the largest, most reliable DDI repository⁸⁴. DDI data were directly downloaded directly from this database.

Type I CC signatures. *Discrete (and discretized) data* (A1–4, B1–5, C1–5, D1, D3, D5, E1–5). Discrete data are expressed as sets of terms, where terms can be proteins, pathways, ATC codes, bit positions of a chemical fingerprint and so on. In some CC spaces, terms are weighted according to their quality or importance (for example, B4 or C5). To convert these sets of terms to a vector form, we applied a protocol originally developed for the numerical representation and comparison of text documents. First, we removed infrequent and frequent terms; that is, terms occurring in less than five and more than m of the molecules ($m = 80\%$ for A1–3, $m = 90\%$ for A4 and $m = 25\%$ for the rest). We then applied a term frequency–inverse document frequency (TF–IDF) transformation to the terms of each molecule, so that ‘term frequency’ was proportional to the weight of the term (when applicable; one otherwise), and the ‘document frequency’ corresponded to the occurrence of the term along the corpus of molecules. As a result of the TF–IDF transformation, less informative terms (that is, promiscuous targets, generic biological processes and so on) become less important. Finally, we applied LSI to the TF–IDF-transformed corpus. LSI is a dimensionality reduction technique based on singular-value decomposition, hence it has parallels with the more popular PCA. In particular, LSI components are also orthogonal and sorted by their contribution to explaining the ‘variance’ of the data. For each dataset, we kept the number of LSI components that explains 90% of the variance. The resulting signatures are thus comprehensive. We also kept track of the ‘elbow’ point of the variance-explained curve (that is, the point of maximum curvature in the scree plot), as this point gives a good trade-off between accuracy (high dimensions) and interpretability (low dimensions).

Continuous data (A5, D2, D4). Data of this type were first robustly scaled column-wise (median, 0 and median absolute deviation, 1; capped at ± 10). Then, for each CC space, we performed a PCA and chose the number of components that explained 90% of the variance. The elbow point was also kept.

Type II CC signatures. We built 25 similarity networks (nodes: molecules, edges: similarities (empirical $-\log_{10}P$ value)). We kept only similarities below a significance $P = 0.01$, and, for each node, we considered at maximum 100 links to other nodes. We ensured that each node was connected to at least three other nodes (ranked by similarity).

We then ran node2vec⁸⁵ to obtain embeddings for each node (molecule) in each network. Node2vec was run with default parameters, that is, $p = 1$, $q = 1$, $k = 10$ (context size), $r = 10$ (walks per source) and $l = 80$ (length of walk). We found an embedding dimension of 128 to be a robust choice across CC spaces (Supplementary Fig. 7).

Clustering and 2D projections of CC signatures. *Clustering.* We used a product-quantized version of k -means⁸⁶ (using product-quantized-table lookups, product-quantized-encoders of 256 bits and eight vector splits) to cluster molecules based on signature similarities. The k -means algorithm requires that a number of clusters k is predefined. We ran k -means with k in the range $2 < k < \sqrt{N}$, N being the number of molecules in the CC space. Inertia (sum of sample distances to centroids) and concentration (inverse of dispersion; that is, number of centroids with another centroid at a significantly close distance ($P < 0.05$)) were calculated at each k . Inertia and concentration curves were smoothed with the Hanning method (window length of $\sqrt{N}/10$) and scaled between zero and one within the explored k range. We chose a k that maximized the geometric mean of both curves, weighting the dispersion curve by the length of the signature with respect to $\sqrt{N}/2$.

2D projections. To have 2D projections of comparable granularity across CC spaces, we performed a k -means clustering on spaces with more than 1,000 samples and took a k of $N/2$, capped at 15,000. Then projections shown in the CCweb and figures of the paper were performed with type I signatures (we observed very similar results with type II signatures). Signatures were projected in a 2D plane using the Barnes–Hut t-SNE algorithm⁸⁷ with a perplexity of 30 and an angle of 0.5. HDBSCAN⁸⁸ was used to identify sparse points (outliers) in the projection. After removing these points, t-SNE was rerun. When necessary, samples were assigned the coordinates of their corresponding centroid.

Correlation between CC spaces. To measure the correlation between two CC spaces, we checked whether, according to the respective CC signatures, molecules in the first space are also similar in the second. We designed a composite correlation coefficient (κ) that quantifies the agreement between the two ranked-similarity lists in several ways (Supplementary Fig. 5). The κ coefficient includes a canonical correlation analysis as well, based on the analysis of dataset cross-covariance and the identification of maximally correlated linear combinations of the two signatures. Thus, high κ values indicate that two CC datasets share similar molecule pairs and that ‘common directions’ can be found between signatures of two datasets.

Given two CC spaces X and Y (paired rows, with $m \times p$ and $m \times q$ dimensions, respectively, where m is the number of common molecules in both CC spaces and p and q are their signature lengths), we did a canonical correlation analysis to identify canonical variables (that is, linear combinations of signature components) that optimally correlate. Correlation between datasets was measured by averaging the Pearson’s correlation between the two first components identified.

We measured the rank-biased overlap (RBO)⁸⁹ between two sorted similarity lists. RBO simulates the behavior of a user scrolling down a list of search results in the web. Higher probabilities of ‘visiting’ a search result are given to higher similarity scores. Two CC spaces with similar RBO lists are thus correlated.

For discretized similarity, when comparing CC spaces pairwise, we classified similarities in the P value intervals $<1 \times 10^{-5}$ (that is, ~ 0), 0.001, 0.01, 0.1, 0.25 and >0.25 . Pairs at these intervals were counted in an ordinal contingency table. Counts were L1-normalized row-wise and column-wise iteratively, and a kappa correlation score was measured on the contingency table using the standard quadratic weighting.

Likewise, we calculated conditional probabilities of two molecules being similar in one CC space when a similarity is observed in another space. The area under the cumulative conditional probabilities (\log_2 -scaled) can be used as a measure of correlation between two spaces.

The correlation measures explained above were unified to a single dataset correlation measure (κ) by simply taking the median value of the individual correlation measures. Values of individual correlations were quantile-normalized before this computation.

Label assignment based on similarity searches. We downloaded drug annotation data from the Drug Repurposing Hub¹⁹ (March 2019). We mapped 5,880 drugs to the CC. These were related to 24 ‘Disease Areas’, 664 ‘Indications’, 1,067 ‘Mechanisms of Action’ and 2,249 ‘Targets’. We then devised the following ‘label assignment’ exercise. For each molecule, we looked for the similar molecules in the dataset using each of the 25 CC spaces separately. Similar molecules were defined as having an empirical $P < 0.001$, calculated on a background specific to the Drug Repurposing Hub. At least three (and at most ten) neighbors (similar molecules) were considered per molecule. Then we evaluated the enrichment (Fisher’s exact test, $P < 0.001$) of labels among the neighbors of the molecules⁹⁰. We required a label to be represented at least in five molecules in the dataset and, at least, in three of the neighbors of the molecule.

For each CC space, we performed independent label assignment exercises for all molecules with known labels. Precision and recall were evaluated, both for CC spaces individually and in a cumulative manner by aggregating (appending) the predictions of each CC space sequentially, spaces being sorted by individual precision.

CC web resource. The CCweb resource (<https://chemicalchecker.org>) is a tool to explore the bioactivity of small molecules, focusing mainly on the identification of compound similarities. To keep pace with its source data, the CCweb will be

updated every 6 months. A simplified representation of the website can be found in Fig. 6. The source code of the resource pipeline is available from the CCweb page.

Home page. The main page of the CCweb consists of a 5×5 grid of panels displaying 2D projections of the signatures (A1–E5). The distribution of all molecules in each CC space is shown as a gray density plot. The user can click on a panel to amplify it. Small-molecule counts and a short explanation of the selected CC space accompany the plot.

Molecules can be queried by InChIKey, PubChem Compound ID (CID) or name. If found in the CC, the compound is shown in the 2D panels where data are available for it. On the right side of the page, a ‘molecule card’ gives basic information about the compound of interest (molecular weight and formula, rule-of-five violations, chemical beauty, popularity, singularity and so on). Of note, a list of targets is given. This list is not meant to be comprehensive, and we encourage users to visit dedicated databases such as ChEMBL to learn more about the targets of their molecules of interest. Targets are sorted by species (human first), then by source ($B1 > B4 > B2 > B5$) and potency (in the case of B4), and finally in alphabetical order.

Regarding libraries and landmark molecules, to facilitate navigation of the 2D panels, we offer the possibility to overlay molecules from the popular chemical collections discussed in this article (APD, PWCK Library and so on; see Fig. 3). These collections can be chosen with the ‘change’ button on the left of the screen. We have selected 100 ‘landmark’ molecules from each collection, since in most cases displaying the full library would be impractical. The selection of the 100 landmark molecules is done such that they are present in as many panels as possible, and favoring their distribution in the 2D projections. To achieve this, we start with the molecule with the highest popularity score. Molecules are sequentially added by bagging always from the most ‘orphan’ CC spaces (that is, the datasets that have the fewest molecules selected). From these, we consider only molecules from the most ‘orphan’ clusters and, among the remaining candidates, we select the one with the highest popularity (that is, more spaces available).

In summary, in the home page of the CCweb, the user can query small molecules and obtain an overview of their location inside the CC. The user will learn the CC spaces where these molecules have data available, with gray 2D density plots indicating whether they are ‘peripheral’ (low-density regions) or ‘central’ (high-density regions). To have a better sense of the location of query molecules, landmark compounds from popular collections can be displayed. Deeper insights can be obtained by clicking on the ‘explore’ button for a molecule of choice.

Explore page. For the list of similar molecules, when a molecule is ‘explored’, we look for similar molecules in the CC database and display them in a 25-column table, corresponding to the CC spaces. In CC spaces where the molecule is available, we measure similarities to other molecules in the space. If the molecule is not available, we infer similarities only against molecules that are present in the space (absent versus present). Inference of similarities is done by simple probabilistic rules using the naive Bayes formulation; that is, we calculate the conditional probabilities of being ‘similar’ ($P < 10^{-5}$, <0.001 , <0.01 , <0.05 , <0.25 and >0.25), based on ‘observed’ similarities in other spaces. The naive Bayes formula was modified by weights⁹¹ to correct for the correlation between spaces and down-weight the individual contribution of strongly correlated CC spaces.

In the ‘explore’ page, measured similarities are shown as filled circles, and inferred similarities as empty ones. Significant similarities ($P < 0.05$) are shown by large colored circles (filled or empty, correspondingly). Small gray circles are shown otherwise for nonsignificant similarities. The list of ‘similar’ molecules can be ranked on the basis of the five levels of complexity (A–E) by clicking on the level name. By default, we up-rank molecules that are similar to the query molecule in many CC spaces, favoring measured similarities over inferred ones. In the CCweb, we give only the top 125 similar molecules (ensuring that at least 25 similar molecules are selected from each of the five CC complexity levels). The user can fetch the full list of similar molecules by clicking on the ‘download’ button on the left of the page.

For the libraries, by default, similar molecules are searched across the CC (all bioactive molecules). The user can choose to search only in certain chemical collections (APD, LINCS and so on). Please note that, in contrast to the main page, we explore the complete collection, not only 100 landmark molecules.

Statistics and help pages. The user can find summary statistics of the CC, plotted as a slideshow in the statistics page. There is also a help page with a short explanation of the resource and a few frequently asked questions. Links to these pages are placed in the black footer of the home page.

Downloads and RESTful access. CC signatures can be downloaded in HDF5 format (‘download’ link in the home page footer). We also provide programmatic access to our signatures through a representational state transfer application programming interface. By default, we provide type II signatures. Type I signatures are available on request.

Compound collections. We downloaded the following chemical collections from ZINC⁹² (<http://zinc15.docking.org>, January 2018): approved drugs (dbap), experimental (dbex) and investigational (dbin) drugs, human metabolites

(hmbendo), traditional Chinese medicines (tcnmp), LINCS compounds (lincs), PWCK Library (pwck), NIH clinical collection (nihcc), NCI diversity collection (ncidiv), and tool compounds (tools). SMILES strings were converted to InChIKeys using the standardization procedure. When an InChIKey was not explicitly present in the CC, we attempted to match the connectivity layer (that is, first 14 characters of the key). Unmatched molecules were discarded. Experimental and investigational drugs were merged into the 'experimental drugs' collection, excluding compounds present in the APD set.

Transcriptional reversion of familial Alzheimer's disease mutations in SH-SY5Y cells. Computational screening.

For the proof of principle, correlation between cancer cell line sensitivity and gene expression reversion. We downloaded cell-sensitivity data from the GDSC²⁵ (GDSC1000 version). We mapped 96 of the GDSC drugs on the gene expression dataset (D1) of the CC. Basal gene expression levels of cancer cell lines were converted to *z* scores based on gene expression values across the panel¹⁹. We took the top 250 over- and underexpressed genes for each cell line, according to the expression *z* score. To measure 'reversion' the direction (up/down) of these two gene sets was flipped.

Cancer cell line transcriptional signatures were then converted to the CC D1 format using the same procedure as that applied to drug signatures; that is, a two-way GSEA of the signatures was done against Touchstone signatures, results were aggregated over Touchstone cell lines and a type I signature was eventually obtained. The signature reversion potential of drugs was calculated by simply measuring the similarity between type I signatures. Signatures having a Pearson's correlation of >0.1 between them and their reversed (flipped) version were excluded from the analysis, that is, they were considered to map poorly to the transcriptional landscape of the CC.

For mutated-versus-WT gene expression signatures, Alzheimer's disease-specific differential gene expression signatures were obtained by comparing the basal gene expression profiles of APP/PSEN1 mutated with WT SH-SY5Y cells (see below). We generated up/down-regulated gene sets conservatively (adjusted $P < 0.01$, \log_2 -fold change $\geq \pm 1.5$) and more permissively ($P < 0.01$, $t \geq \pm 2$). Additional versions of the signatures were obtained by keeping only genes related to Alzheimer's disease and Tau pathology in OpenTargets (confidence scores of 0.5 (high), 0.2 (medium) and 0.1 (low)). Finally, composite signatures were also derived by simply measuring intersection or union of gene sets (for example, consensus PSEN1^{M146V} signatures could be obtained from the homozygous and heterozygous clones; see below) (Supplementary Data 1).

For signature reversion, all the signatures above were flipped and converted to the D1 format. Reversion potential (connectivity) of CC compounds was then measured by similarity of type I signatures to each Alzheimer's disease-related signature. Connectivity scores were robustly normalized (median and median absolute deviation), and aggregated when necessary with the tertile statistic as indicated in ref. ²². Full results are given in Supplementary Data 1.

CRISPR-Cas9 gene edition for Alzheimer's disease cell models. Single-guide RNA (sgRNA) sequences targeting APP and PSEN1 were designed using the Zhang laboratory CRISPR design tool (<http://crispr.mit.edu>), and cloned into a modified version of pX330 plasmid expressing green fluorescent protein and puromycin resistance⁹⁴. Next, 200 long single stranded donor oligonucleotides (ssODN) were used as a template for inducing HDR and designed to introduce the desired mutation together with silent mutations to protect both the ssODN template and also the mutated allele once homologous recombination has taken place (Supplementary Fig. 10a). ssODN were purchased from ITD with phosphodiester modification in the 3'. All sequences are listed in Supplementary Table 3.

SH-SY5Y cells were cultured in DMEM/F12 (1:1) medium supplemented with 10% fetal bovine serum (FBS), glutamine and antibiotics (Thermo Fisher Scientific). For transfection, the SH-SY5Y cells were seeded in T-75 flasks and allowed to grow to 80% confluency. A mixture of X330 plasmid and ssODN template was transfected using linear polyethylenimine (Polysciences) in Opti-MEM medium (Thermo Fisher Scientific) supplemented with 10% FBS. Three days after transfection, cells were trypsinized and seeded again in the presence of 2 $\mu\text{g ml}^{-1}$ puromycin (Sigma-Aldrich). Selection pressure with puromycin was kept for 1 week, and then selected cells were allowed to expand and recover for 1–2 weeks. Some of the cells were then used to measure overall HDR efficiency and the rest were single-cell cloned in 96-well plates using a FACSaria II flow cytometer (BD Biosciences). Three to four weeks after cloning, confluent wells were split into two 96-well plates, one to expand the clone and the other to analyze the genotype. DNA extraction was performed adding 50 μl of DirectPCR-tail lysis reagent (VWR) supplemented with 0.4 mg ml^{-1} of proteinase K (Roche), and plates were incubated overnight at 55 °C. The next day, lysates were moved to 96-well PCR plates and we inactivated Proteinase K incubating at 85 °C for 40 min. Next, 5 μl of lysate was used to amplify by PCR the genomic region surrounding the edition target with recombinant Taq DNA polymerase (Thermo Fisher Scientific), followed by digestion with restriction enzymes to screen for the introduction of new restriction sites encoded in the ssODN template (Supplementary Fig. 10). Using either genomic DNA (gDNA primers) or reverse transcribed RNA (complementary DNA primers) as template, mutated cells were routinely tested for the presence of the mutation by selective digestion with restriction enzymes

(Supplementary Fig. 10b) of PCR-amplified DNA fragments surrounding the mutation. All primer sequences are listed in Supplementary Table 3. All restriction enzymes were purchased from New England Biolabs.

Of the clones isolated, we obtained homozygous clones (APP^{V717F/V717F} and PSEN1^{M146V/M146V}) and also heterozygous mutants in which the second allele had a two-nucleotide deletion leading to a displacement in the reading frame and a premature stop codon, therefore called "null" (APP^{V717F/null}) or, in the case of the M146V mutation, a three-nucleotide deletion encoding for an amino acid deletion at position 149 (PSEN1^{M146V/L149Δ}). We then measured the two main forms of A β peptide secretion (A β 42 and A β 40) in all the isolated clones, and we observed an increase in the A β 42/A β 40 ratio (Supplementary Fig. 10c). All cell lines are available on request from the authors.

A β quantification. A β peptides were quantified by enzyme-linked immunosorbent assay- (ELISA)-based assays using either the 6E10 A β Triplex by MesoScale Diagnostics or the Wako ELISA kit Human β Amyloid (1–40) and Wako ELISA kit Human β Amyloid (1–42) High-Sensitivity, following the manufacturer's instructions. Direct comparison of the results showed similar results for the two quantification assays.

Drug treatment of SH-SY5Y clones. Cells were differentiated for 3 d in neurobasal medium supplemented with B27, glutamax (all Thermo Fisher Scientific), 10 μM retinoic acid (Sigma-Aldrich) and 50 ng ml^{-1} brain-derived neurotrophic factor (Peprotech). Then medium was renewed in the presence of the indicated concentration of drugs. All drugs were dissolved in DMSO, and controls of cells treated with DMSO were run in parallel, at a final concentration of 0.1% DMSO. After 3 d, supernatants were stored at –80 °C for A β measurement and cells were either incubated for 1 h in the presence of 3-(4,5-dimethylthiazol-2-yl)-2,5-diphenyltetrazolium bromide and lysed in DMSO to check the viability, or lysed with RLT buffer (Qiagen) for RNA extraction using the RNeasy mini kit (Qiagen). Three independent experiments were performed.

Gene expression. To obtain the signature profile derived from these FAD mutated cells, SH-SY5Y WT and mutated cells were differentiated for 6–7 d in the presence of retinoic acid and Brain-derived neurotrophic factor to recapitulate a phenotype more similar to neurons as previously described²⁶. The secretion of A β followed the same pattern as that observed in nondifferentiated cells. Samples of purified RNA of WT (APP/PSEN1^{WT/WT}), APP^{V717F/null}, PSEN1^{M146V/L149Δ} and PSEN1^{M146V/M146V} (clone no. 2) were extracted and submitted to the IRB Functional Genomics Facility, where sample quality was assessed using an Agilent Bioanalyzer. The whole-genome expression profile was generated using Affymetrix PrimeView arrays. Three independent experiments were used to obtain the expression profiles. All gene expression signatures have been deposited in the GEO ([GSE137202](https://www.ncbi.nlm.nih.gov/geo/query/acc.cgi?acc=GSE137202)).

Evaluation of results. The 'reversion capacity' of tested drugs was measured as follows. We ranked the differential gene expression results of the treated-versus-untreated comparison performed on mutated cells (ranked list 1). In parallel, we ranked the differential gene expression results of the mutated-versus-WT comparison (ranked list 2). In Fig. 4b, we simply traverse ranked list 2 (*x* axis, from both tails) and measure the identification of genes at the other end (top 250) of ranked list 1 (*y* axis). Ranked list 1 was randomized to assign significance to observations. To obtain a 'reversion strength' value per gene, we designed a reversion score based on the difference in ranks between mutated-versus-WT (reference) and treated-versus-control gene expression profiles. We scaled reversion scores, ranging from –1 (underexpressed genes in the reference are overexpressed on treatment) to +1 (vice versa). To test whether reversed genes were enriched in Alzheimer's disease genes (OpenTargets score >0.5), we performed a weighted Kolmogorov–Smirnov, taking as weights the absolute value of the reversion score⁹⁵.

Identification of small-molecule mimetics of biologics against IL2R, IL-12 and EGFR. Computational screening. When screening biodrug-related signatures, biodrugs were defined by their targets (that is, IL2R, IL12B and EGFR). We derived D1 (transcriptional), C3 (pathway), C4 (biological process) and C5 (interactome) CC signatures for these three targets. C3 and C4 signatures were obtained by simply mapping pathways and biological processes of the targets, respectively, and expressing them as CC signatures by TF-IDF/LSI transformation. Regarding C5 signatures, we applied the same procedure than the one applied to compounds; that is, we mapped target neighbors in interactomes using HotNet2 and then we obtained the corresponding type I signature. For D1, we downloaded gene expression signatures from shRNA experiments obtained from LINCS L1000, and mapped them analogously to small-molecule perturbations.

For matching biological-related signatures, we devised a computational screening for drugs in D1 spaces and in any of the C3–5 spaces. We asked for candidates to be among the top 250 drugs in terms of similarity of the target signature to at least one of the C3–5 spaces, and ranked them on the basis of similarity of transcriptional profiles (that is, mimicking D1 similarity).

Cells. PBMC were purchased from StemCell Technologies and maintained in RPMI medium supplemented with 10% FBS, glutamine and antibiotics (Thermo Fisher

Scientific). Prestimulated PBMC were obtained by culturing 10^6 PBMCs per ml with $0.5 \mu\text{g ml}^{-1}$ soluble anti-CD28 (CD28.2) and anti-CD3 (OKT3) antibodies (Thermo Fisher Scientific) for 3 d. After this stimulation period, cells were washed and left untreated for 3 d before restimulation. H1650, Jurkat and MT-4 cells were cultured in RPMI supplemented with 10% FBS, glutamine and antibiotics (Thermo Fisher Scientific). A431 and HeLa cells were cultured in DMEM supplemented with 10% FBS, glutamine and antibiotics (Thermo Fisher Scientific). NK-92 cells were purchased from ATCC (CRL-2407), cultured in alpha-MEM without ribo- and deoxyribo-nucleosides (Thermo Fisher Scientific), and supplemented with FBS and horse serum (both Thermo Fisher Scientific), 0.2 mM inositol (Sigma-Aldrich), 0.1 mM 2-mercaptoethanol (Sigma-Aldrich) and 0.02 mM folic acid (Sigma-Aldrich), penicillin/streptomycin and glutamine (both Thermo Fisher Scientific). Every 2–3 d, 100 U ml^{-1} of recombinant IL-2 (Peprotech) was added.

PBMC proliferation assay. Resting and prestimulated PBMC were loaded with $2 \mu\text{M}$ CFSE (Thermo Fisher Scientific) in PBS with 0.1% FBS for 7 min at 37°C . After two washes in complete medium, PBMC were pretreated for 1 h with the corresponding drugs, followed by stimulation with 0.5 ng ml^{-1} IL-2 or $5 \mu\text{g ml}^{-1}$ PHA (Sigma-Aldrich). Three days after stimulation, cell fluorescence was measured using a Gallios Flow Cytometer (Beckton Coulter). Analysis was conducted with FlowJo software.

Phospho-STAT5 quantification by flow cytometry. Prestimulated PBMC pretreated for 1 h with the corresponding compounds were stimulated for 20 min with 0.5 ng ml^{-1} IL-2, fixed with Fix Buffer I (BD Biosciences), permeabilized with Perm Buffer III (BD Biosciences), and finally stained with a PE-labeled antiphospho-Stat5 (pY694, BD Biosciences). Staining was measured in a Gallios Flow Cytometer and analysis was conducted with FlowJo software. Two compounds, SU11652 and Z55175877 showed autofluorescence at high concentrations when cells were analyzed by flow cytometry, therefore STAT5 phosphorylation was measured by western blot as detailed next.

Proliferation of cell lines. Here, 5×10^4 Jurkat or MT-4 cells were incubated in the presence of the indicated compounds. In the case of HeLa cells, 5×10^3 cells were seeded the day before the compounds were added to the supernatant. After 3 d, cells were incubated for 1 h in the presence of 3-(4,5-dimethylthiazol-2-yl)-2,5-diphenyltetrazolium bromide and viability/proliferation was quantified as indicated above.

IL-12 stimulation. The day before the experiment, NK-92 cells were counted, washed with RPMI supplemented with 10% FBS and seeded in 24-well plates (3×10^5 cells per well) in RPMI 10% FBS in the absence of IL-2. On the day of the experiment, cells were pre-incubated with the indicated compounds for 1 h and then stimulated with 50 ng ml^{-1} IL-12 (Peprotech). Cells were pelleted after 1 h of stimulation and lysed for western blot analysis or kept up to 5 h in culture. They were then pelleted and RNA was extracted for quantitative PCR (qPCR) analysis as indicated next.

qPCR. Purified RNA samples were reverse transcribed with the High-Capacity cDNA Reverse Transcription Kit (Thermo Fisher Scientific) and qPCR was performed in a QuantStudio 6 Flex Real-Time PCR System (Thermo Fisher Scientific) using the LightCycler 480 SYBR Green I Master mix (Roche). *Ct* values were normalized using *GAPDH* as a reference gene and the $\Delta\Delta C_t$ method to quantify the fold change of the gene of interest. Primers are shown in Supplementary Table 3.

EGFR analysis. Here, 0.15×10^6 A431 or H1650 cells were seeded in 24-well plates the day before the experiment. Cells were treated with the indicated inhibitors for 24 h. Cells were then washed with PBS and lysed for western blot analysis.

Western blot. Cells stimulated with or without cytokines were washed in PBS, concentrated and resuspended in lysis buffer (50 mM Tris-HCl (pH 7.5), 1 mM EGTA, 1 mM EDTA, 1% (wt/wt) Triton X-100) supplemented with protease inhibitor cocktail (Roche) and phosphatase inhibitor cocktail (Roche). Lysates were subjected to SDS-PAGE in Mini-PROTEAN TGX Stain-Free Precast Gels (Biorad) and transferred to a polyvinylidene difluoride membrane using the Trans-Blot Turbo Transfer System (Biorad). Images of developed blots were acquired with the Chemidoc Touch Imaging System (Biorad).

The following antibodies were used for immunoblotting: horseradish peroxidase-conjugated secondary antibodies (Thermo Fisher Scientific), anti-Actin (Merck), anti-Stat5 (D206Y, Cell Signaling Technology) and anti-Stat4 (C46B10, Cell Signaling Technology). Phosphospecific antibodies recognizing phospho-Tyr693 of Stat4 and phospho-Tyr694 of Stat5 (D47E7) were also from Cell Signaling Technology. The EGFR monoclonal antibody (clone 13) was purchased from BD Biosciences.

Statistical analysis. Data were analyzed with the Prism statistical package. Unless otherwise indicated in the figure legend, *P* values were calculated using an unpaired, one-tailed Student *t*-test.

Reporting Summary. Further information on research design is available in the Nature Research Reporting Summary linked to this article.

Data availability

All gene expression signatures have been deposited in the GEO ([GSE137202](https://www.ncbi.nlm.nih.gov/geo/query/acc.cgi?acc=GSE137202)).

Code availability

To facilitate access to our data, we built a web-based resource (<https://chemicalchecker.org>), which includes all the bioactivity signatures in HDF5 format and the full code of the CC resource.

References

- Axen, S. D. et al. A simple representation of three-dimensional molecular structure. *J. Med. Chem.* **60**, 7393–7409 (2017).
- Bemis, G. W. & Murcko, M. A. The properties of known drugs. 1. Molecular frameworks. *J. Med. Chem.* **39**, 2887–2893 (1996).
- Durant, J. L., Leland, B. A., Henry, D. R. & Nourse, J. G. Reoptimization of MDL keys for use in drug discovery. *J. Chem. Inf. Comput. Sci.* **42**, 1273–1280 (2002).
- Lipinski, C. A. Lead- and drug-like compounds: the rule-of-five revolution. *Drug Discov. Today Technol.* **1**, 337–341 (2004).
- Congreve, M., Carr, R., Murray, C. & Jhoti, H. A ‘rule of three’ for fragment-based lead discovery? *Drug Discov. Today* **8**, 876–877 (2003).
- Wishart, D. S. et al. DrugBank 5.0: a major update to the DrugBank database for 2018. *Nucleic Acids Res.* **46**, D1074–D1082 (2018).
- Cheng, H. et al. ECoD: an evolutionary classification of protein domains. *PLoS Comput. Biol.* **10**, e1003926 (2014).
- Gilson, M. K. et al. BindingDB in 2015: a public database for medicinal chemistry, computational chemistry and systems pharmacology. *Nucleic Acids Res.* **44**, D1045–D1053 (2016).
- Hastings, J. et al. ChEBI in 2016: improved services and an expanding collection of metabolites. *Nucleic Acids Res.* **44**, D1214–D1219 (2016).
- Thiele, I. et al. A community-driven global reconstruction of human metabolism. *Nat. Biotechnol.* **31**, 419–425 (2013).
- Cerami, E. G. et al. Pathway Commons, a web resource for biological pathway data. *Nucleic Acids Res.* **39**, D685–D690 (2011).
- Fabregat, A. et al. The Reactome Pathway Knowledgebase. *Nucleic Acids Res.* **46**, D649–D655 (2018).
- Pryszcz, L. P., Huerta-Cepas, J. & Gabaldon, T. MetaPhOrs: orthology and paralogy predictions from multiple phylogenetic evidence using a consistency-based confidence score. *Nucleic Acids Res.* **39**, e32 (2011).
- Kruger, F. A. & Overington, J. P. Global analysis of small molecule binding to related protein targets. *PLoS Comput. Biol.* **8**, e1002333 (2012).
- Zwierzyńska, M. & Overington, J. P. Classification and analysis of a large collection of in vivo bioassay descriptions. *PLoS Comput. Biol.* **13**, e1005641 (2017).
- Szklarczyk, D. et al. The STRING database in 2017: quality-controlled protein–protein association networks, made broadly accessible. *Nucleic Acids Res.* **45**, D362–D368 (2017).
- Li, T. et al. A scored human protein–protein interaction network to catalyze genomic interpretation. *Nat. Methods* **14**, 61–64 (2017).
- Kanehisa, M., Sato, Y., Kawashima, M., Furumichi, M. & Tanabe, M. KEGG as a reference resource for gene and protein annotation. *Nucleic Acids Res.* **44**, D457–D462 (2016).
- Kandasamy, K. et al. NetPath: a public resource of curated signal transduction pathways. *Genome Biol.* **11**, R3 (2010).
- Mi, H. et al. PANTHER version 11: expanded annotation data from Gene Ontology and Reactome pathways, and data analysis tool enhancements. *Nucleic Acids Res.* **45**, D183–D189 (2017).
- Kelder, T. et al. WikiPathways: building research communities on biological pathways. *Nucleic Acids Res.* **40**, D1301–D1307 (2012).
- Mosca, R., Ceol, A. & Aloy, P. Interactome3D: adding structural details to protein networks. *Nat. Methods* **10**, 47–53 (2013).
- Leiserson, M. D. et al. Pan-cancer network analysis identifies combinations of rare somatic mutations across pathways and protein complexes. *Nat. Genet.* **47**, 106–114 (2015).
- Iorio, F. et al. Discovery of drug mode of action and drug repositioning from transcriptional responses. *Proc. Natl. Acad. Sci. USA* **107**, 14621–14626 (2010).
- Barretina, J. et al. The Cancer Cell Line Encyclopedia enables predictive modelling of anticancer drug sensitivity. *Nature* **483**, 603–607 (2012).
- Basu, A. et al. An interactive resource to identify cancer genetic and lineage dependencies targeted by small molecules. *Cell* **154**, 1151–1161 (2013).
- Chabner, B. A. NCI-60 cell line screening: a radical departure in its time. *J. Natl. Cancer Inst.* **108**, djv388 (2016).
- Azur, M. J., Stuart, E. A., Frangakis, C. & Leaf, P. J. Multiple imputation by chained equations: what is it and how does it work? *Int. J. Meth. Psychiatr. Res.* **20**, 40–49 (2011).

76. Nelson, J. et al. MOSAIC: a chemical-genetic interaction data repository and web resource for exploring chemical modes of action. *Bioinformatics* **34**, 1251–1252 (2017).
77. Wawer, M. J. et al. Toward performance-diverse small-molecule libraries for cell-based phenotypic screening using multiplexed high-dimensional profiling. *Proc. Natl Acad. Sci. USA* **111**, 10911–10916 (2014).
78. Brown, A. S. & Patel, C. J. A standard database for drug repositioning. *Sci. Data* **4**, 170029 (2017).
79. Piñero, J. et al. DisGeNET: a comprehensive platform integrating information on human disease-associated genes and variants. *Nucleic Acids Res.* **45**, D833–D839 (2017).
80. Kuhn, M., Letunic, I., Jensen, L. J. & Bork, P. The SIDER database of drugs and side effects. *Nucleic Acids Res.* **44**, D1075–1079 (2016).
81. Kuhn, M. et al. Systematic identification of proteins that elicit drug side effects. *Mol. Syst. Biol.* **9**, 663 (2013).
82. Duran-Frigola, M. & Aloy, P. Analysis of chemical and biological features yields mechanistic insights into drug side effects. *Chem. Biol.* **20**, 594–603 (2013).
83. Davis, A. P. et al. The Comparative Toxicogenomics Database: update 2017. *Nucleic Acids Res.* **45**, D972–D978 (2017).
84. Ryu, J. Y., Kim, H. W. & Lee, S. Y. Deep learning improves prediction of drug–drug and drug–food interactions. *Proc. Natl Acad. Sci. USA* **115**, 4304–4311 (2018).
85. Grover, A. & Leskovec, J. node2vec: scalable feature learning for networks. Preprint at <https://arxiv.org/abs/1607.00653> (2016).
86. Matsui, Y. O., Yamasaki, K. & Aizawa, T. K PQk-means: billion-scale clustering for product-quantized codes. Preprint at <https://arxiv.org/abs/1709.03708> (2017).
87. Maaten, L. v. d. Barnes–Hut-SNE. Preprint at <https://arxiv.org/abs/1301.3342> (2013).
88. McInnes, L. & Healy, J. Accelerated hierarchical density based clustering. *Proc. 2017 IEEE International Conference on Data Mining Workshops* (IEEE, 2017).
89. Webber, W., Moffat, A. & Zobel, J. A similarity measure for indefinite rankings. *ACM Trans. Inf. Syst.* **28**, 1–38 (2010).
90. Lo, Y. C. et al. Large-scale chemical similarity networks for target profiling of compounds identified in cell-based chemical screens. *PLoS Comput. Biol.* **11**, e1004153 (2015).
91. Rennie, J. D. M., Shih, L., Teevan, J. & Karger, D. R. Tackling the poor assumptions of naive Bayes text classifiers. *Proc. International Conference on Machine Learning* 616–623 (AAAI Press, 2003).
92. Irwin, J. J. & Shoichet, B. K. ZINC—a free database of commercially available compounds for virtual screening. *J. Chem. Inf. Model* **45**, 177–182 (2005).
93. Fernandez-Torras, A., Duran-Frigola, M. & Aloy, P. Encircling the regions of the pharmacogenomic landscape that determine drug response. *Genome Med.* **11**, 17 (2019).
94. Badia, R. et al. SAMHD1 is active in cycling cells permissive to HIV-1 infection. *Antiviral Res.* **142**, 123–135 (2017).
95. Saxena, V., Orgill, D. & Kohane, I. Absolute enrichment: gene set enrichment analysis for homeostatic systems. *Nucleic Acids Res.* **34**, e151 (2006).

Acknowledgements

We thank the SB&NB laboratory members for their support and helpful discussions. We are grateful to the Broad Institute and National Center for Advancing Translational Sciences (NCATS-NIH) for providing compounds on request, and J. Duran-Frigola for the website design. We also thank the IRB Barcelona Biostatistics and Bioinformatics Unit and the IRB Functional Genomics Facility. P.A. acknowledges the support of the Spanish Ministerio de Economía y Competitividad (grant no. BIO2016-77038-R), the INB/ELIXIR-ES (grant no. PT17/0009/0007), the European Research Council (SysPharmAD, grant no. 614944) and ‘La Caixa’ BioMedTec (grant no. CTEC_15).

Author contributions

M.D.-F., E.P. and P.A. designed the study, analyzed the results and wrote the manuscript. M.D.-F. did the computational analysis, together with M.B., T.J.-B. and D.A. O.G.-P. implemented the web server. E.P. and V.A. carried out the experimental validations. All authors have read and approved the manuscript.

Competing interests

The authors declare no competing interests.

Additional information

Supplementary information is available for this paper at <https://doi.org/10.1038/s41587-020-0502-7>.

Correspondence and requests for materials should be addressed to M.D.-F. or P.A.

Reprints and permissions information is available at www.nature.com/reprints.

Reporting Summary

Nature Research wishes to improve the reproducibility of the work that we publish. This form provides structure for consistency and transparency in reporting. For further information on Nature Research policies, see [Authors & Referees](#) and the [Editorial Policy Checklist](#).

Statistics

For all statistical analyses, confirm that the following items are present in the figure legend, table legend, main text, or Methods section.

n/a Confirmed

- ☐ ☒ The exact sample size (n) for each experimental group/condition, given as a discrete number and unit of measurement
- ☐ ☒ A statement on whether measurements were taken from distinct samples or whether the same sample was measured repeatedly
- ☐ ☒ The statistical test(s) used AND whether they are one- or two-sided
Only common tests should be described solely by name; describe more complex techniques in the Methods section.
- ☒ ☐ A description of all covariates tested
- ☐ ☒ A description of any assumptions or corrections, such as tests of normality and adjustment for multiple comparisons
- ☐ ☒ A full description of the statistical parameters including central tendency (e.g. means) or other basic estimates (e.g. regression coefficient) AND variation (e.g. standard deviation) or associated estimates of uncertainty (e.g. confidence intervals)
- ☐ ☒ For null hypothesis testing, the test statistic (e.g. F , t , r) with confidence intervals, effect sizes, degrees of freedom and P value noted
Give P values as exact values whenever suitable.
- ☒ ☐ For Bayesian analysis, information on the choice of priors and Markov chain Monte Carlo settings
- ☒ ☐ For hierarchical and complex designs, identification of the appropriate level for tests and full reporting of outcomes
- ☒ ☐ Estimates of effect sizes (e.g. Cohen's d , Pearson's r), indicating how they were calculated

Our web collection on [statistics for biologists](#) contains articles on many of the points above.

Software and code

Policy information about [availability of computer code](#)

Data collection

qPCR data was collected using QuantStudio 6 Flex Real-Time PCR System. Arrays were scanned with GeneChip scanner GSC3000 (Affymetrix). Western blot images were acquired using Chemidoc Touch Imaging System (Biorad).

Data analysis

Flow cytometry was analyzed using FlowJo v10. Figures and statistical analysis were made using Graphpad Prism 7. Heat propagation of the C3-5 spaces was done with HotNet 2.0.

For manuscripts utilizing custom algorithms or software that are central to the research but not yet described in published literature, software must be made available to editors/reviewers. We strongly encourage code deposition in a community repository (e.g. GitHub). See the Nature Research [guidelines for submitting code & software](#) for further information.

Data

Policy information about [availability of data](#)

All manuscripts must include a [data availability statement](#). This statement should provide the following information, where applicable:

- Accession codes, unique identifiers, or web links for publicly available datasets
- A list of figures that have associated raw data
- A description of any restrictions on data availability

The Chemical Checker code is available at http://gitlabssnb.irsbbarcelona.org/project-specific-repositories/chemical_checker. The data produced by our pipeline can be accessed from <http://chemicalchecker.org>. All gene expression signatures have been deposited in GEO (GSE137202).

Field-specific reporting

Please select the one below that is the best fit for your research. If you are not sure, read the appropriate sections before making your selection.

☒ Life sciences ☐ Behavioural & social sciences ☐ Ecological, evolutionary & environmental sciences

For a reference copy of the document with all sections, see [nature.com/documents/nr-reporting-summary-flat.pdf](https://www.nature.com/documents/nr-reporting-summary-flat.pdf)

Life sciences study design

All studies must disclose on these points even when the disclosure is negative.

| | |
|-----------------|--|
| Sample size | No sample-size calculation was performed. Sample sizes were determined by consistency of measurable differences. |
| Data exclusions | No data were excluded. |
| Replication | Several independent experiments were performed to ensure reproducibility, and all the replicas were significantly correlated. |
| Randomization | No formal randomization was done. When required (e.g microarray experiments), samples were distributed in groups to prevent batch effects affecting our conclusions. |
| Blinding | Researchers performing the RNAseq experiments were blinded during the processing of experimental samples. |

Reporting for specific materials, systems and methods

We require information from authors about some types of materials, experimental systems and methods used in many studies. Here, indicate whether each material, system or method listed is relevant to your study. If you are not sure if a list item applies to your research, read the appropriate section before selecting a response.

Materials & experimental systems

| n/a | Involved in the study |
|-------------------------------------|---|
| <input type="checkbox"/> | <input checked="" type="checkbox"/> Antibodies |
| <input type="checkbox"/> | <input checked="" type="checkbox"/> Eukaryotic cell lines |
| <input checked="" type="checkbox"/> | <input type="checkbox"/> Palaeontology |
| <input checked="" type="checkbox"/> | <input type="checkbox"/> Animals and other organisms |
| <input checked="" type="checkbox"/> | <input type="checkbox"/> Human research participants |
| <input checked="" type="checkbox"/> | <input type="checkbox"/> Clinical data |

Methods

| n/a | Involved in the study |
|-------------------------------------|--|
| <input checked="" type="checkbox"/> | <input type="checkbox"/> ChIP-seq |
| <input type="checkbox"/> | <input checked="" type="checkbox"/> Flow cytometry |
| <input checked="" type="checkbox"/> | <input type="checkbox"/> MRI-based neuroimaging |

Antibodies

| | |
|-----------------|---|
| Antibodies used | Anti-mouse (goat) IgG HRP conjugate (Thermo Fisher Scientific, ref G-21040, 1:5,000); Anti-rabbit (goat) IgG HRP conjugated (Thermo Fisher Scientific, ref 65-6120, 1:5,000); anti-Actin (Clone C4; Merck; ref MAB1501; Lot 3018859; 1:20,000); anti-Stat5 (Clone D206Y; Cell Signaling Technology; ref 94205S; Lot 3; 1:1,000); anti-Stat4 (Clone C46B10; Cell Signaling Technology; ref 2653S; Lot 3; 1:1,000); anti-pStat4(Y693) (Clone D2E4; Cell Signaling Technology; ref 4134S; Lot 4; 1:1,000); anti-pStat5(Y694) (Clone D47E7; Cell Signaling Technology; ref 4322S; Lot 8; 1:1,000); Anti-EGFR monoclonal antibody (Clone 13; BD Biosciences; ref 610017; Lot 4020981; 1:1,000); anti-CD28 (Clone CD28.2; Thermo Fisher Scientific; ref 16-0289-81; Lot 1928813; 0.5 ug/ml); anti-CD3 (clone OKT3; Thermo Fisher Scientific; ref 16-0037-81; Lots 1981083/1952722; 0.5 ug/ml); PE-labelled anti pStat5(pY694) (Clone47; BD Biosciences; 562077 ; Lot 7278790; 1:20). Anti-CD25 (Daclizumab; Novus Biologicals; ref NBP2-52660; Lot T1713A15; different concentrations indicated in the figures). |
| Validation | Cell Signalling Technologies provides examples of internal validation for western blot application of the antibodies (https://www.cellsignal.com/): anti-Stat5 (Clone D206Y; ref 94205S), anti-Stat4 (Clone C46B10, ref 2653S), anti-pStat4(Y693) (Clone D2E4; ref 4134S) and anti-pStat5(Y694) (Clone D47E7; ref 4322S). Merck-Millipore provides the quality control for western blot application for each of their anti-Actin lots (ref MAB1501) in www.merckmillipore.com . BD Biosciences routinely tests the anti-EGFR monoclonal antibody (ref 61001) for western blot application and the PE-labelled anti pStat5(pY694) (56207) for flow cytometry applications, as stated in www.bdbiosciences.com . Thermo Fisher Scientific (https://www.thermofisher.com) provides example figures for the anti-CD28 (ref 16-0289-81) and the anti-CD3 (ref 16-0037-81) antibodies, both for flow cytometry applications and functional assays. They also provide a list of references that have previously used these antibodies. For all the cases, our results at western blot level and flow cytometry level matched with the ones expected according to the results provided by the supplier. |

Eukaryotic cell lines

Policy information about [cell lines](#)

| | |
|--|--|
| Cell line source(s) | NK-92 and HeLa (ATCC); Jurkat and MT-4 cells (José A Esté, Irsicaixa Institute); A431 (Ernest Giralt, IRB Barcelona); H1650 (Frank Supek, IRB Barcelona); SH-SY5Y (Jens Lüders, IRB Barcelona) |
| Authentication | Cells were not authenticated. |
| Mycoplasma contamination | Mycoplasma tests were performed routinely for all the cell lines and all were negative. |
| Commonly misidentified lines (See ICLAC register) | No commonly misidentified cell lines were used in the study. |

Flow Cytometry

Plots

Confirm that:

- ☒ The axis labels state the marker and fluorochrome used (e.g. CD4-FITC).
- ☒ The axis scales are clearly visible. Include numbers along axes only for bottom left plot of group (a 'group' is an analysis of identical markers).
- ☒ All plots are contour plots with outliers or pseudocolor plots.
- ☒ A numerical value for number of cells or percentage (with statistics) is provided.

Methodology

| | |
|---------------------------|--|
| Sample preparation | Proliferation of PBMC was quantified by direct acquisition of CFSE pre-loaded PBMC. For intracellular stainings, PBMC were fixed with Fix Buffer I (BD Biosciences), permeabilized with Perm Buffer III (BD Biosciences), and finally stained with a PE-labelled anti-phospho-Stat5 (pY694; BD Biosciences) following manufacturer instructions. |
| Instrument | Gallios Flow Cytometer (Beckton Coulter) |
| Software | FlowJo v10 |
| Cell population abundance | N/A |
| Gating strategy | N/A |

☐ Tick this box to confirm that a figure exemplifying the gating strategy is provided in the Supplementary Information.