1

2

3

4          **CDEK: Clinical Drug Experience Knowledgebase**

5

6

7     Rebekah H. Griesenauer[1], Constantino Schillebeeckx[1], Michael S. Kinch[1*]

8

9

10

11    [1] Center for Research Innovation in Biotechnology, Washington University in St. Louis, Missouri,

12    United States

13

14

15    * Corresponding author

16    E-mail: michael.kinch@wustl.edu (MK)

17

18

19

20

21

22

23

24

25

**ABSTRACT**

The Clinical Drug Experience Knowledgebase (CDEK) is a database and web platform of active pharmaceutical ingredients with evidence of clinical testing as well as the organizations involved in their research and development. CDEK was curated by disambiguating intervention and organization names from ClinicalTrials.gov and cross-referencing these entries with other prominent drug databases. Approximately 43% of active pharmaceutical ingredients in the CDEK database were sourced from ClinicalTrials.gov and cannot be found in any other prominent compound-oriented database. The contents of CDEK are structured around three pillars: active pharmaceutical ingredients (n = 22,292), clinical trials (n = 127,223), and organizations (n = 24,728). The envisioned use of the CDEK is to support the investigation of many aspects of drug development, including discovery, repurposing opportunities, chemo- and bio-informatics, clinical and translational research, and regulatory sciences.

**Database URL: http://cdek.wustl.edu**

**INTRODUCTION**

The process in which drugs are discovered and developed has fundamentally changed since the inception of the pharmaceutical industry and continues to evolve. Several research groups have peered into the past to identify trends in pharmaceutical innovation based upon FDA approved medicines (1–3). The Center for Research Innovation in Biotechnology (CRIB) at Washington University in St. Louis is amidst an ongoing effort to objectively track and analyze trends in the innovation of new medicines. Several published works were facilitated by analysis of a precursor database (curated and maintained by CRIB) of all FDA-approved new molecular entities (NMEs), which included their mechanistic basis, therapeutic applications, and organizations guiding their clinical development. This NME database (http://cribdb.wustl.edu) also includes products that were

50   once approved but no longer marketed as a result of toxicity, lack of efficacy, obsolescence,

51   production issues, or lack of demand.

52

53   A handful of reviews on the biopharmaceutical industry trends and innovation sources revealed a

54   trove of findings, many unexpected, and all supported by objective data (all of which we have made

55   public). As one example, a handful of organizations have recently come to control two-thirds of NMEs

56   and these marketing organizations often have little or no internal drug discovery or development

57   activities (4). Whereas large, traditional pharmaceutical companies receive most FDA approvals,

58   upstart biotechnology companies increasingly dominate early-stage discovery (including patents

59   and Investigational New Drug (IND) applications) (5). The NME database also revealed the causes

60   and impact of corporate consolidation in transforming research and development. Whereas 60% of

61   all acquired biotechnology companies were acquired within 5 years (before or after) their first NME

62   approval was granted, the number of new organizations to receive their first approval has not kept

63   pace (6). Consequently, the net number of research organizations that remain active and

64   independent in new drug research has has eroded from over 200 firms in 2004 to 100 firms at the

65   end of 2015 (7).

66

67   Based on findings with FDA-approved medicines, we analyzed the mechanistic basis and therapeutic

68   indications of FDA approved medicines and changes over time. In some cases, these works

69   emphasized therapeutic areas (e.g., the decline in anti-infectives or the rise in oncology (8)) while

70   others focused upon drug targets, revealing three target families dominate FDA-approved drugs (G-

71   protein coupled receptors, membrane channels and transporters, and targets involving nuclear

72   signaling (9)). Beyond clinical indications and drug targets, exploration of other facets of the

73   biotechnology industry enabled by the NME database included regulatory pathways and timelines

74   (10), vaccine development (11), and the rise of biologics (12).

3

75  Although intriguing, we considered prior observations of pharmaceutical research and development

76  trends to be undoubtedly skewed by focusing only upon FDA-approved medicines. It is generally

77  understood most drug research does not conclude with a single FDA-approval as post-approval

78  research (e.g., additional indications or post-approval commitments) capture an ever-increasing

79  fraction of research and development expenditures and are not captured in analyses of drugs based

80  solely upon a designation of "FDA-approved."  Compounding the problem, the timelines required for

81  drug development mean an FDA-approval reflects research and development activities that were

82  likely initiated more than a decade before, enfeebling any analyses intended to assess current or

83  predict future research and development activity. Consequently, conjectures and definitive

84  conclusions are not feasible absent a more comprehensive accounting of drug development efforts;

85  including an assessment of successes, failures, and those experimental medicines currently being

86  developed.

87

88  Powerful insights can be obtained by analyzing and modeling drug "failures." In Gayvert et al. (13), a

89  random forest machine learning algorithm classified a set of compounds as "FDA approved" or

90  "Failed for Toxicity" based on chemical structure and drug target features. In this study, 784 FDA

91  approved drugs and 100 "toxic" drugs were used to train and validate the machine learning model.

92  Ideally, failed drugs would have made up a higher percentage of the sample, but sufficient data on

93  failed drugs are not readily available. Nonetheless, these findings revealed machine learning

94  predictions can be quite powerful provided that they are supplied with enough data for training and

95  validation. Wong et al. (14) were able to assign a probability of success to clinical trials solely by

96  following drugs through clinical trial phase transitions and comparing intended medical applications.

97  The data for this study was limited to information from a commercial dataset and not available

98  publicly. While an open assessment of all experimental medicines would be preferable, the authors

99    stated "trained analysts would require tens of thousands of hours of labor" (14) to perform such a

100   study using ClinicalTrials.gov, a public source for clinical trials data.

101

102   The current lack of public data on successful, failed, and on-going drug studies sparked the

103   development of the Clinical Drug Experience Knowledgebase (CDEK- http://cdek.wustl.edu) with the

104   purpose of creating a public platform to analyze all active pharmaceutical ingredients that have ever

105   been tested in humans, as well as their sponsoring organizations and those participating in pre-

106   approval clinical activites. Based on insights derived from previous studies, we focused on three

107   primary pillars for the first instantiation of CDEK: active pharmacetucial ingredients, organizations,

108   and clinical trials. Each pillar is shown in Figure 1 with surrounding metadata fields. Foreign keys in

109   the database link each pillar together. In the next section, we review the current state of clinical stage

110   pharmaceuticals available in public databases.

> **Figure 1:** Overview of CDEK contents with three primary pillars: Active Pharmaceutical Ingredients,
> Organizations, and Clinical Trials. Each metatopic is surrounded with the current fields (solid lines) and
> planned metadata fields (dashed lines).

111

### *Current state of clinical stage pharmaceuticals in public databases*

113   Several biopharmaceutical databases have emerged over the last decade to enable chemo- and bio-

114   informatics research in the field of drug discovery, including chemical structures to support *in silico*

115   drug discovery, drug repurposing opportunities, and trends in the drug development enterprise. A

116   decade ago, fewer than 200 peer-reviewed articles were published per year referencing a

117   biopharmaceutical database. Today, over 2,500 articles annually cite biopharmaceutical databases

118   and this rate continues to grow exponentially. We recently surveyed several open and freely available

119   databases to explore the current landscape of clinical stage pharmaceuticals and found a collection

120   of databases having drug records that display some evidence of clinical experience.

121

122    A selection of databases is listed in Table 1, including a brief description of the clinical content of the

123    database. However, these databases often contain discovery-level or preclinical molecules that have

124    never or will ever enter the clinic.  The PubChem (15) database, housing over 100 million compound

125    records, can be filtered to clinical stage compounds by extracting records sourced from

126    ClinicalTrials.gov, ToxCast, or the NCATS Pharmaceutical Collection. ChEMBL (16), another large

127    compound database, can be filtered to clinical stage compounds by selecting records with a

128    max_phase greater or equal to one (with max_phase corresponding to the farthest clinical trial phase

129    the compound has been registered). DrugBank (17), an encyclopedia of active pharmaceutical

130    ingredients, can be filtered to clinical compounds by selecting "Approved", "Withdrawn",

131    "Investigational", "Illicit", or "Nutraceutical" from their "Drug Group" metadata field. Other databases

132    focus explicitly on approved or withdrawn medicines, making their whole catalog of drugs relevant

133    in terms of clinical experience.

134

135    In a study that inspired the creation of CDEK, our group downloaded the clinical-stage active

136    pharmaceutical ingredients from the sources listed in Table 1. Approximately 11,760 unique active

137    pharmaceutical ingredients with evidence of clinical experience were available collectively from

138    those data sources.  However, the total number of active pharmaceutical ingredients that have ever

139    been tested in humans was likely much higher. For example, Wong *et al*. used the Informa Pharma

140    Intelligence databases "TrialTrove" and "Pharmaprojects" to complete their study on estimating

141    clinical trial success rates. In their study, they cited extracting over 21,143 unique compounds from

142    the Informa Pharma Intelligence databases with corresponding clinical trial information (14).

143

144    Table 1: Public databases containing clinical stage active pharmaceutical ingredients.

| Database | Scope | Clinical Experience Evidence | Access |
|---|---|---|---|
| **PubChem** | Chemical entities and their bioactivities | Records sourced from Clinicaltrials.gov, ToxCast, or NCATS Pharmaceutical Collection | https://pubchem.ncbi.nlm.nih.gov |
| **ChEMBL** | Bioactivity for drug discovery | Field "max_phase" =>1 | https://www.ebi.ac.uk/chembl |
| **DrugBank** | *in silico* drug discovery and exploration | Field "DRUG GROUP" = "Approved OR Withdrawn OR Investigational OR Illicit OR Nutraceutical" | https://www.drugbank.ca |
| **DrugCentral** | Active pharmaceutical ingredients approved by FDA and other agencies | All records are approved or withdrawn medicines | http://drugcentral.org |
| **SuperDrug2** | Marketed drugs | All records are approved or withdrawn medicines | http://cheminfo.charite.de/superdrug2 |
| **CRIB NME** | FDA approved molecular entities and biopharmaceutical organizations | All records are approved or withdrawn medicines | http://cribdb.wustl.edu |
| **repoDB** | Drug repurposing | All records are either approved or have been in clinical trials. | http://apps.chiragjpgroup.org/repoDB |
| **WITHDRAWN** | Withdrawn or discontinued drugs | All records are withdrawn medicines | http://cheminfo.charite.de/withdrawn |

145

146     Such findings suggest other active pharmaceutical ingredients may exist in the public domain but

147     have not been curated. ClinicalTrials.gov (accessed through the Aggregate Analysis of Clinical Trials

148     (AACT) database), for example, contains over 286,811 unique trials with over 246,005 unique

149     "intervention names" in a trial (as of 10/20/2018). Multiple "intervention names" correspond to the

150     same active pharmaceutical ingredient. To achieve the ambitious goal of "studying all drugs ever

151     tested in a human", it was necessary to mine and disambiguate ClinicalTrials.gov data to supplement

152     the compounds available in current open access drug databases.

153

154     Descriptions of the disambiguation of ClinicalTrials.gov interventions and organizations follow.

155     Detail on how other databases were used to cross-reference unique ClinicalTrials.gov interventions

156     is also summarized. CDEK is the culmination of this curation effort and is a public database and web

157     platform to interrogate all active pharmaceutical ingredients where there exists objective evidence

158     of human clinical testing. CDEK aggregates metadata surrounding active pharmaceutical ingredients,

159     including the details of clinical trial design, intended indications, and organizations responsible for

160     development. The envisioned use of the CDEK is to support the investigation of many aspects of drug

7

161    development, including discovery, repurposing opportunities, chemo- and bio-informatics, clinical

162    and translational research, and regulatory sciences. The platform is intended to serve a wide

163    audience interested in investigational agents, which have reached clinical stage development. The

164    uses enabled by CDEK also include the elucidation of broad or focused trends, competitive

165    intelligence, improving drug development efficiency and conveying best practices of lessons learned

166    and future directions.

167

168    **METHODS**

169    ***CDEK Construction: Curating ClinicalTrials.gov data***

170    Construction of CDEK arose from multiple iterations beginning with the predominant source of data:

171    ClinicalTrials.gov accessed through the Aggregate Analysis of Clinical Trials (AACT) database (18).

172    ClinicalTrials.gov is a repository of clinical trial registrations in the United States and is maintained

173    by the National Library of Medicine (NLM) at the National Institutes of Health (NIH) in collaboration

174    with the Food and Drug Administration (FDA). The AACT database was developed and is maintained

175    by the Clinical Trials Transformation Initiative (CTTI) group, a government-academic collaboration

176    between the FDA and Duke University. The AACT database contains ClinicalTrials.gov data that has

177    been parsed and deposited into a structured relational database. AACT also links clinical trials data

178    to Medical Subject Headings (MeSH terms), a controlled vocabulary containing terms describing

179    disease indications and interventions. This mapping enables querying the data by intervention and

180    disease indication terms. In this first step, we were primarily interested in removing the ambiguity

181    in the trial intervention names and names of sponsoring organizations.

182

183    The AACT *interventions* table has the field *intervention_type* with the following distinct terms used to

184    describe an intervention in a trial: Drug, Behavioral, Diagnostic Test, Dietary Supplement, Other,

185    Device, Biological, Procedure, Combination Product, Genetic, and Radiation. To initially populate

8

186    CDEK with therapeutic clinical trials, all AACT pharmaceutical interventions were included whereas

187    interventions labeled Behavioral, Diagnostic Test, Device, Radiation or Other were excluded. CDEK

188    was populated with associated clinical trial data and organizations linked to those entries. The

189    organizations in turn were parsed from the *sponsors* table*, overall_officials* table, and

190    *responsible_parties* table within AACT. Collectively, these tables contain the lead and collaborating

191    sponsors, trial affiliation data for various study roles (e.g. Principal Investigator, Study Chair), and

192    trial affiliation data for the party type (e.g. Sponsor, Sponsor-Investigator).

193

194    In a first round of data cleanup, the names of active pharmaceutical ingredients and organizations

195    were validated. Each active pharmaceutical ingredient was manually labeled by biomedical research

196    curators as being one of either *Vaccine, Gene therapy, Cell therapy, Small molecule, Biologic*

197    *(synthesized in organisms or cell lines), Biological (derived from human material), Animal product* or

198    *Botanical; and* any active pharmaceutical ingredient not categorized as such was removed from the

199    dataset. Additionally, active pharmaceutical ingredient names were manually curated and any active

200    pharmaceutical ingredient listed as a combination drug was split into its constituent parts. Manual

201    validation and cleaning of active pharmaceutical ingredient names included correcting obvious

202    misspellings and removing salt or solvent forms. Similarly, each organization was labeled as being

203    one of *Individual*, *Academic/Hospital, Government, Foundation, For profit* or *Unknown,* and each

204    organization name was validated and normalized to have consistent naming nomenclature. Figure 2

205    illustrates an example of the curation process for an active pharmaceutical ingredient.

206

207

208

**Figure 2:** An example that illustrates the process of extracting inteventions from ClinicalTrials.gov (through AACT) and creating a unique active pharmaceutical ingredient record in CDEK. Curation begins by extracting the intervention names from trials containing active pharmaceutical ingredients and cleaning names to strip any perfulous text (e.g. dosing amount, dosing freqency). Once complete, an automated program flags entities that should be merged into a single CDEK record using a set of "merging" criteria. The curation software will also flag entities that are made up of two or more active pharmaceutical ingredients using a set of "splitting" criteria (e.g. the drug "Mavyret" is a combination of two active pharmaceutical ingredients, glecaprevir and pibrentasvir, used to treat hepatitis C). A unique CDEK active pharmaceutical ingredient record is created and assigned a unique id, a type, and a preferred name. All names are stored as synonymns and all trials are linked to the unique active pharmaceutical ingredient ID. Finally, several external databases are cross-referenced to pull metadata and provide hyperlinks to more information about that active pharmaceutical ingredient. This metadata was also used to flag entries that should be merged into a single active pharmaceutical ingredient.

209

210 *Construction: Cross-referencing with public biopharmaceutical databases*

211 Additional sources of data were ingested into the database following the first round of cleanup.

212 Several open drug-compound databases containing clinically tested therapeutics to capture active

213 pharmaceutical ingredients with evidence of clinical testing outside of the ClinicalTrials.gov registry.

214 These databases included Drugbank(17), ChEMBL(16), PubChem(15), SuperDrug2(19),

215 DrugCentral(20), WITHDRAWN(21), repoDB (22) and CRIB NME (4). The first three of these

216 databases were subsetted to access only those therapeutics with evidence of clinical testing, while

217 the remainder contain soley clinically-tested therapeutics (approved by a regulatory agency,

218 withdrawn from the market for any reason, or associated with a clinical trial). All DrugBank (v5.0.7)

219 compounds labeled "experimental" were excluded from CDEK as DrugBank defines "experimental"

220 as "drugs that are at the preclinical or animal testing stage." The ChEMBL database labels drug

1|

221 compound records as having a *max_phase*, the maximum clinical trial phase for which that drug

222 compound has been tested. Any compounds with a *max_phase* greater than 0 was ingested from

223 ChEMBL v23. Any PubChem compound annotated as sourced from ClinicalTrials.gov were ingested.

224 Additionally, all approved drugs listed on the regulatory websites (as of April 2018) of the Food and

225 Drug Administration (Drugs@FDA) and European Medicines Agency (EMA) were parsed, validated

226 and ingested. The metadata provided by these external databases were used to facilitate the

227 disambiguation process described in the next section.

228

229 ***Construction: Removing ambiguity to get a list of unique Interventions and Organizations***

230 After initial cleanup and ingestion, expert curators split and merged organizations and active

231 pharmaceutical ingredients based on their metadata. We performed this cleanup and ingestion

232 process semi-manually by first programatically flagging data for review followed by manual

233 validation of each flagged entry. The program identified active pharmaceutical ingredients to be

234 considered for merging when two or more distinct entries are were labelled with the same active

235 pharmaceutical ingredient name, *source_api_id* (the ID given to the active pharmaceutical ingredient

236 in a given source), chemical structure (SMILES string), or had overlapping synonyms. Similarly, the

237 program flagged records for splitting active pharmaceutical ingredients into multiple distinct

238 compounds when multiple non-distinct chemical structure data was associated with a given active

239 pharmaceutical ingredient or if multiple *source_api_id*s were associated with the active

240 pharmaceutical ingredient. The program calculated similarity scores (e.g. Levenshtein distance) for

241 all pairs of organizations to identify highly similar organizations pairs, which expert curators then

242 manually validated as either being the same organization or not.

243

244 Figure 3 demonstrates an example of the ambiguous nature of ClincalTrials.gov data. Our particular

245 home institution, Washington University in St. Louis (WUSTL), was designated by more than 50

1

246    unique representations in ClinicalTrials.gov. This represents the ambiguity challenge to be remedied.

247    Figure 3 shows a network in which all red nodes are different representations of the WUSTL name

248    and all black nodes are the clinical trials associated with that name. After disambiguation, all WUSTL

249    affiliated trials were represented as one organization: "Washington University in St. Louis". The June

250    2017 snapshot of AACT has 54047 organization names associated with the 127,220 clinical trials in

251    CDEK. We manually validated  and collapsed these entries into 24,728 unique CDEK organizations.

252    Furthermore, AACT has 104,627 unique interventions names that we manually validated and

253    collapsed to 17,096 CDEK active pharmaceutical ingredients. During the curation process, we stored

254    all names, which had been collapsed into single organizations as "alternative names". This allows for

255    users to search many different terms in our web application.

**Figure 3:** Network graph of trials associated with Washington University in St. Louis. The left graph

shows different representations of Washington University in St. Louis in ClinicalTrials.gov as red nodes.

Examples of different names representing "Washington University in St. Louis" include: "Washington

University School of Medicine", "Washington Universite Siteman Cancer Center", and various misspellings

of the word 'university'. Black nodes are the clinical trials associated with each different name for the

Washington University in St. Louis organization. The right graph shows CDEK data with Washington

University in St. Louis as a single organization with its corresponding clinical trials.

256

257    *CDEK Contents*

258    Table 2 provides summary statistics of CDEK contents: active pharmaceutical ingredients (n =

259    22,292), clinical trials (n = 127,223), and organizations (n = 24,728).

260

261    Table 2: Summary counts of CDEK data

| Organization Type | Count | API Type | Count | Trial Phase | Count |
|---|---|---|---|---|---|
| Academic/Hospital | 9495 | Small Molecules | 13169 | Phase 2 | 32538 |
| For Profit | 6577 | Biologics | 2583 | Phase 1 | 23656 |
| Individual | 3634 | Botanicals | 1769 | Phase 3 | 22641 |
| Unknown | 3183 | Vaccine | 1698 | N/A | 18830 |
| Foundation | 1200 | Cell Therapy | 1521 | Phase 4 | 18267 |
| Government | 658 | Biological | 1182 | Phase 1/Phase 2 | 7054 |
| **Total Orgs** | **24747** | Animal Product | 233 | Phases 2/Phase 3 | 3184 |
| | | Gene Therapy | 157 | Early Phase 1 | 1163 |
| | | **Total APIs** | **22312** | **Total Trials** | **127333** |

262

263    CDEK includes all prophylactic and therapeutic chemical or biological entities, including but not

264    limited to vaccines, cell therapies, gene therapies, animal products, and biologics – many of which are

265    not typically included in other popular compound-oriented databases.

266

267    **RESULTS**

268    ***CDEK Platform***

269    The CDEK platform used the open-source web framework, Django, which follows the model-view-

270    controller architectural pattern. This allows the internal representation of data (the models) to be

271    separated from the presentation to the end user (the view). In the back-end, the models were

272    implemented as a PostgreSQL database and all data is hosted on Heroku. The controller and views

273    rendered the front-end of the platform using a mix of HTML, CSS and javascript.

274

275    The CDEK platform provided two query functionalities, allowing users to quickly interface with the

276    data without having any prior familiarity with a structured query language (SQL). The first

277    functionality, a *basic search* (http://cdek.wustl.edu/search/) enable the user to do a fuzzy, case-

1

278     insensitive search for keywords or synonyms in order to find either active pharmaceutical

279     ingredients or organizations. This functionality serves as a quick, simplified means of interacting with

280     a single datum. The result displays summary statistics of the basic CDEK pillars. Figure 4 shows an

281     example of active pharmaceutical ingredients and organization summary pages. For an active

282     pharmaceutical ingredient, the clinical trial distribution is plotted according to trial phase and

283     organizations involved in its developed is plotted according to organization type. For an organization,

284     the involvement in clinical trials and active pharmaceutical ingredient development is plotted

285     according to trial phase and active pharmaceutical ingredient type, respectively. In both search

286     displays, a list of alternative names is given. For those interested in the source data, or who seek to

287     visualize the ingested reference, CDEK allows the user to link to external cross-referencing databases.

**Figure 4**: Example Active Pharmaceutical Ingredient and Organization summary pages from the CDEK platform. Adalimumab, was the top selling drug of 2017 while GlaxoSmithKline has the most associated clinical trials in the CDEK database.

288     Users are directed to an advanced query functionality to access the granular CDEK data.

289

290     The *advanced query* functionality (cdek.wustl.edu/query/) provides users with more control over

291     the metadata are used to filter the dataset. A dynamically generated user-interface (UI) allows a user

292     to build a SQL-like query, in a *WISYWIG* ("what you see is what you get") fashion, without having any

293     previous knowledge of SQL. Complex queries can be quickly generated by building filtration rules

294     (predicates) and by combining them with boolean logic. These data are then submitted to the back-

295     end through an AJAX (Asynchronous JavaScript And XML) call to a database-view which combines all

296     the CDEK data into a single table. This AJAX call initializes a Celery worker which will process the

297     query request on a separate Heroku worker dyno and return the result in a non-blocking fashion;

298     this ensures that the platform can scale properly as more queries are submitted and ensures a better

299     user experience. Results are presented in a familiar table-like manner with sortable columns and

300     hyperlinks to individual data instances. A RESTful API (application programming interface) provides

301     an endpoint for viewing these individual data when either requesting a single active pharmaceutical

302     ingredient or organization instance. This endpoint dynamically generates interactive charts which

303     summarize the data for the given data instance. Our advanced query builder allows a user to filter

304     CDEK data to granular details. Figure 5 shows a screen shot of the query tool web application. In this

305     example, the data returned will be all unique Phase III clinical trials (n = 681) studying lung or

306     cardiovascular diseases, excluding vaccines, and run by GlaxoSmithKline as the lead sponsor between

> **Figure 5:** Our advanced query builder allows users to filter down CDEK data to very granular details. In this example, the data returned will be all unique Phase III clinical trials studying lung or cardiovascular diseases, excluding vaccines, that were ran by GlaxoSmithKline as the lead sponsor between 2012 and 2017.

307     2012 and 2017.

308

309     ***Lessons Learned***

310     Approximately 17,096 unique active pharmaceutical ingredients in CDEK were sourced from

311     ClinicalTrials.gov, 9,781 of which currently cannot be found in any databases cross-referenced in

312     CDEK (see Table 1). These active pharmaceutical ingredients comprise 3160 small molecules, 1477

313     vaccines, 1438 cell therapies, 1387 biologics, 1084 botanicals, 982 biologicals, 143 gene therapies,

314     and 110 animal products. The databases included for initial cross-referencing primarily focus on

315     small molecules and biologics. Therefore, we reviewed unique small molecules and biologics

316     extracted from ClinicalTrials.gov, hereafer refered to as "unique CDEK records". Most (90%) unique

317     CDEK records have been registered in three or fewer clinical trials and 85% of the clinical trials

318     referencing these drugs are prior to Phase III. This indicates that early stage active pharmaceutical

319     ingredients might not typically be flagged for curation in traditional databases. Another interesting

320     trend is almost two-thirds (64%) of the unique CDEK records were sponsored by for profit

321   organizations. This contrasts to the whole CDEK dataset where less than one third (30%) of all trial

322   lead sponsors are for profit organizations.

323

324   The active pharmaceutical ingredient contents of CDEK was compared with other common

325   compound-oriented drug databases including: PubChem, Chembl, DrugBank, DrugCentral,

326   SuperDrug2, WITHDRAWN, repoDB, and drugs@FDA. Despite our initial assumption that existing

327   databases, once aggregated, would convey a comprehensive list of experimental medicines,

328   approximately 43% of active pharmaceutical ingredients in the CDEK database were extracted from

329   AACT and cannot be found in any of the other compound-oriented databases listed above.

330

331   We reviewed the overlap of active pharmaceutical ingredients with evidence of clinical testing among

332   several open databases, including those listed in Table 1, AACT, and the Drugs@FDA database. Figure

333   6 shows the this overlap as a heatmap, comparing content across several drug databases. This

334   visualization demonstrates that some databases are almost complete subsets of others (99% of

335   repoDB compounds can be found in ChEMBL, DrugCentral and DrugBank). PubChem, one of the

336   largest compound libraries showed consistently high overlap values across the spectrum. The

337   overlap between AACT active pharmaceutical ingredients and PubChem is the highest, closely

338   followed by AACT and ChEMBL.

> **Figure 6:** Heatmap displaying the overlap in active pharmaceutical ingredients (APIs) between any two
> databases in CDEK. The coloring and number displayed at the intersection between any two databases is
> the total number of shared APIs. The total number of unique APIs from each database that has evidence of
> clinical experience is noted in paranthesis next to each database name label.

339

340   **DISCUSSION**

16

341    The purpose of CDEK is to provide researchers with an open database and platform to study the

342    entire drug development enterprise by interrogating *all* active pharmaceuticals with evidence of

343    clinical testing. While not truly comprehensive, we have created the first release of such a resource

344    and below we discuss several on-going strategies for improvement.

345

346    The first instantiation of CDEK was derived from a June 2017 snapshot of the AACT database. Over

347    20,000 trials registered in ClinicalTrials.gov were not included in the first instantiation but we are

348    currently developing a novel "ingestion pipeline" to allow curators to update the data automatically

349    and in real time. Databases listed as cross-referencing sources will be updated in CDEK in the future

350    along with the addition of new data sources – such as ToxCast and ZINC. Future curated databases

351    will also be merged into CDEK under the conditions they are public, verifiable and contain evidence

352    of clinical-trial candidates.

353

354    The curation of several new metadata fields will be incorporated into CDEK. These fields are

355    summarized in Figure 1 encircled by dashed lines. They include information such as patents

356    surrounding active pharmaceutical ingredients, approval status of each indication associated with an

357    active pharmaceutical ingredient, clinical trial study results, and the merger and acquisition activity

358    of for-profit organizations conducting clinical trials.

359

360    Another on-going area of development is mining scientific publications containing clinical trial

361    information. ClinicalTrials.gov was created in response to the Food and Drug Administration

362    Modernization Act of 1997 (FDAMA), with the first public version of ClinicalTrials.gov released in

363    2000. Therefore, it is necessary to search public reports of clinical studies for trials that may not have

364    been registered, or that were conducted prior to 1997.

365

366  Finally, continued efforts are being made to clean and disambiguate any residual errors propagated

367  through the initial data cleanup. We intend to employ higher standards for chemical data set curation

368  methods, such as those outlined by Fourches et al (23). Due to the expansive efforts needed to keep

369  CDEK up-to-date and accurate, our group is also interested in deploying crowd-based curation

370  methods in the future.

371

372  **CONTACTING CDEK**

373  CDEK was developed and is maintained by the Center for Research Innovation in Biotechnology

374  (CRIB) at Washington University in St. Louis. CRIB studies the blend of science, business, and

375  regulation of biotechnology, medical devices, and healthcare IT to ensure continued improvements

376  in the delivery of medical innovations and public health. CRIB is actively pursuing collaborations to

377  study the data within CDEK. Errors and suggestions for improvement can be submitted at

378  http://cdek.wustl.edu/about/. Or contact us via e-mail at cdek at wustl dot edu.

379

380  **ACKNOWLEDGEMENTS**

387

388  **References**

389  1.   Munos B. Lessons from 60 years of pharmaceutical innovation. Nat Rev Drug Discov.

390       2009;8(12):959–68.

391  2.   Vitaku E, Smith DT, Njardarson JT. Analysis of the structural diversity, substitution patterns,

392       and frequency of nitrogen heterocycles among U.S. FDA approved pharmaceuticals. J Med

393       Chem [Internet]. 2014;57(24):10257–74. Available from:

394       https://www.scopus.com/inward/record.uri?eid=2-s2.0-

395       84920194283&doi=10.1021%2Fjm501100b&partnerID=40&md5=81299f9be3f00c0415a0

396       72a7338b4ad3

397  3.   Wu P, Nielsen TE, Clausen MH. FDA-approved small-molecule kinase inhibitors. Trends

398       Pharmacol Sci [Internet]. 2015;36(7):422–39. Available from:

399       https://www.scopus.com/inward/record.uri?eid=2-s2.0-

400       84933181941&doi=10.1016%2Fj.tips.2015.04.005&partnerID=40&md5=37bcb861c4cd2e1

401       b9fd01c14344923c4

402  4.   Kinch MS, Haynesworth A, Kinch SL, Hoyer D. An overview of FDA-approved new molecular

403       entities: 1827–2013. Drug Discov Today. 2014;19(8):1033–9.

404  5.   Kinch MS. The rise (and decline?) of biotechnology. Drug Discov Today. 2014;19(11):1686–

405       90.

406  6.   Kinch MS. Post-approval fate of pharmaceutical companies. Drug Discov Today.

407       2015;20(2):170–4.

408  7.   Griesenauer RH, Kinch MS. 2016 in review: FDA approvals of new molecular entities. Drug

409       Discov Today. 2017;

410  8.   Kinch MS, Merkel J, Umlauf S. Trends in pharmaceutical targeting of clinical indications:

411       1930–2013. Drug Discov Today. 2014;19(11):1682–5.

412  9.   Kinch MS, Hoyer D, Patridge E, Plummer M. Target selection for FDA-approved medicines.

413       Drug Discov Today. 2015;20(7):784–9.

414  10.  Patridge E V, Gareiss PC, Kinch MS, Hoyer DW. An analysis of original research contributions

415       toward FDA-approved drugs. Drug Discov Today. 2015;20(10):1182–7.

416    11.    Griesenauer RH, Kinch MS. An overview of FDA-approved vaccines & their innovators.

417            Expert Rev Vaccines. 2017;16(12):1253–66.

418    12.    Kinch MS. An overview of FDA-approved biologics medicines. Drug Discov Today. 2015;

419    13.    Gayvert KM, Madhukar NS, Elemento O. A Data-Driven Approach to Predicting Successes and

420            Failures of Clinical Trials. Cell Chem Biol. 2016;23(10):1294–301.

421    14.    Wong CH, Siah KW, Lo AW. Estimation of clinical trial success rates and related parameters.

422            Massachusetts Institute of Technology, Department of Electrical Engineering and Computer

423            Science; 2017.

424    15.    Kim S, Thiessen PA, Bolton EE, Chen J, Fu G, Gindulyte A, et al. PubChem substance and

425            compound databases. Nucleic Acids Res. 2015;44(D1):D1202–13.

426    16.    Gaulton A, Hersey A, Nowotka M, Bento AP, Chambers J, Mendez D, et al. The ChEMBL

427            database in 2017. Nucleic Acids Res. 2016;45(D1):D945–54.

428    17.    Wishart DS, Feunang YD, Guo AC, Lo EJ, Marcu A, Grant JR, et al. DrugBank 5.0: a major

429            update to the DrugBank database for 2018. Nucleic Acids Res. 2017;46(D1):D1074–82.

430    18.    Tasneem A, Aberle L, Ananth H, Chakraborty S, Chiswell K, McCourt BJ, et al. The database for

431            aggregate analysis of ClinicalTrials. gov (AACT) and subsequent regrouping by clinical

432            specialty. PLoS One. 2012;7(3):e33677.

433    19.    Siramshetty VB, Eckert OA, Gohlke B-O, Goede A, Chen Q, Devarakonda P, et al. SuperDRUG2:

434            a one stop resource for approved/marketed drugs. Nucleic Acids Res. 2017;46(D1):D1137–

435            43.

436    20.    Ursu O, Holmes J, Knockel J, Bologa CG, Yang JJ, Mathias SL, et al. DrugCentral: online drug

437            compendium. Nucleic Acids Res. 2016;gkw993.

438    21.    Gillespie LD, Gillespie WJ, Robertson MC, Lamb SE, Cumming RG, Rowe BH. WITHDRAWN:

439            Interventions for preventing falls in elderly people. Cochrane database Syst Rev.

440            2009;(2):CD000340-CD000340.

441    22.    Brown AS, Patel CJ. A standard database for drug repositioning. Sci data. 2017;4:170029.

442    23.    Fourches D, Muratov E, Tropsha A. Trust, but verify: On the importance of chemical structure

443            curation in cheminformatics and QSAR modeling research. Journal of Chemical Information

444            and Modeling. 2010.

445

446

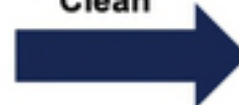Figure

Figure

Representation of Washington University in St. Louis

Associated Clinical Trial

Ambiguity Removed

Figure

# adalimumab Report issue

Biologics

ChEMBL | DrugBank | Drug Central | Drugs@FDA | KEGG | MeSH | repoDB

## Trials (330 total)



## Organizations (257 total)



## Metadata

Alternative names

adalimumab | adalimumab-atto | adalimumab (genetical recombination) | amjevita | d2e7 | humira | lu200134 | lu-200134
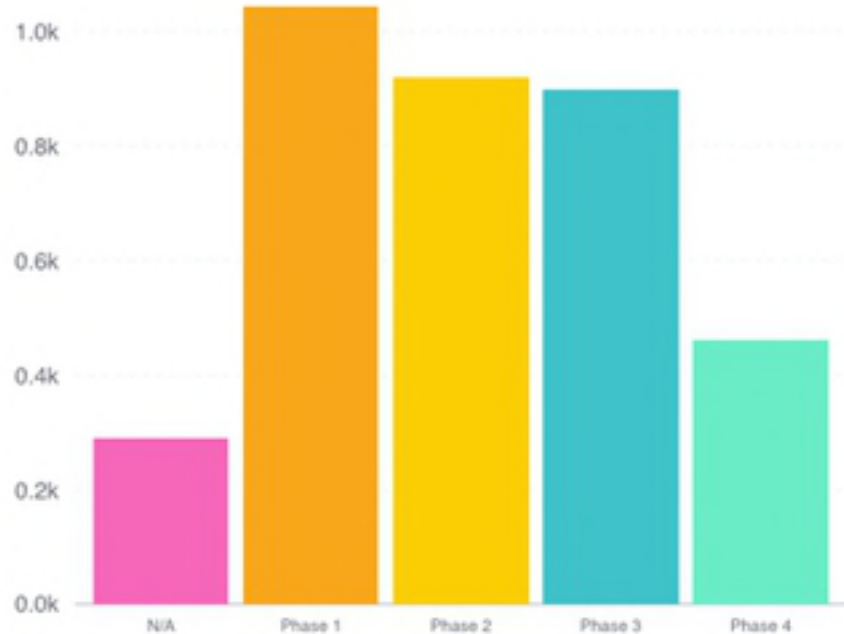
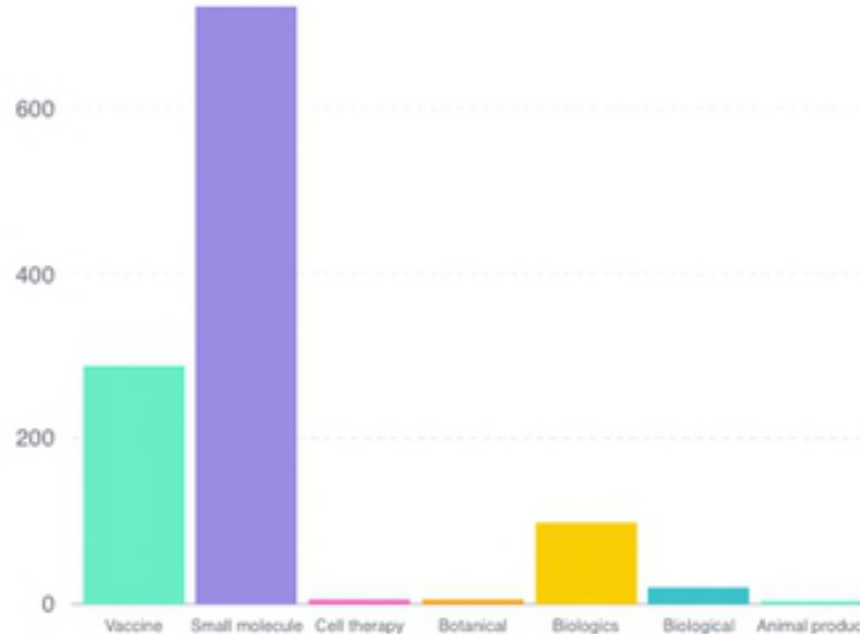# GlaxoSmithKline ✎ Report issue

For profit

**Founded on:** 2000

**Founded in:** London United Kingdom

gsk do more feel better live longer

## Trials (3,612 total)



## APIs (1,148 total)



### Alternative names

Barrier Therapeutics/ Stiefel, a GSK Company | GlasoSmithKline | GlaxoSmithKline AG, Switzerland | Glaxosmithkline Biologicals S.A. | Glaxosmithkline/Quintiles | GSK-CIHR Research Chair in Respiratory Health Care Delivery, | GSK Vaccines Institute for Global Health (GVGH) | Stiefel, a GSK Company

# Figure

Figure

# Counts of APIs common between any two sources

| | ChEMBL (n = 7038) | DrugCentral (n = 4608) | DrugBank (n = 4742) | Drugs@FDA (n = 1845) | PubChem (n = 10023) | WITHDRAWN (n = 618) | CRIB NME (n = 1951) | repoDB (n = 1541) | SuperDrug2 (n = 3911) |
|---|---|---|---|---|---|---|---|---|---|
| AACT (n = 17096) | 4279 | 2324 | 3716 | 1436 | 6200 | 251 | 1508 | 1355 | 1860 |
| ChEMBL (n = 7038) | | 4282 | 3692 | 1804 | 5985 | 592 | 1756 | 1529 | 3643 |
| DrugCentral (n = 4608) | | | 2681 | 1695 | 4398 | 582 | 1645 | 1532 | 3824 |
| DrugBank (n = 4742) | | | | 1642 | 4029 | 509 | 1691 | 1524 | 2281 |
| Drugs@FDA (n = 1845) | | | | | 1679 | 362 | 1583 | 1254 | 1359 |
| PubChem (n = 10023) | | | | | | 614 | 1651 | 1453 | 3808 |
| WITHDRAWN (n = 618) | | | | | | | 374 | 265 | 573 |
| CRIB NME (n = 1951) | | | | | | | | 1297 | 1400 |
| repoDB (n = 1541) | | | | | | | | | 1315 |

Figure