

SciCura

Shared, General-purpose Curation
of Life Science Knowledge



"A bold new paradigm for how research information is disseminated"

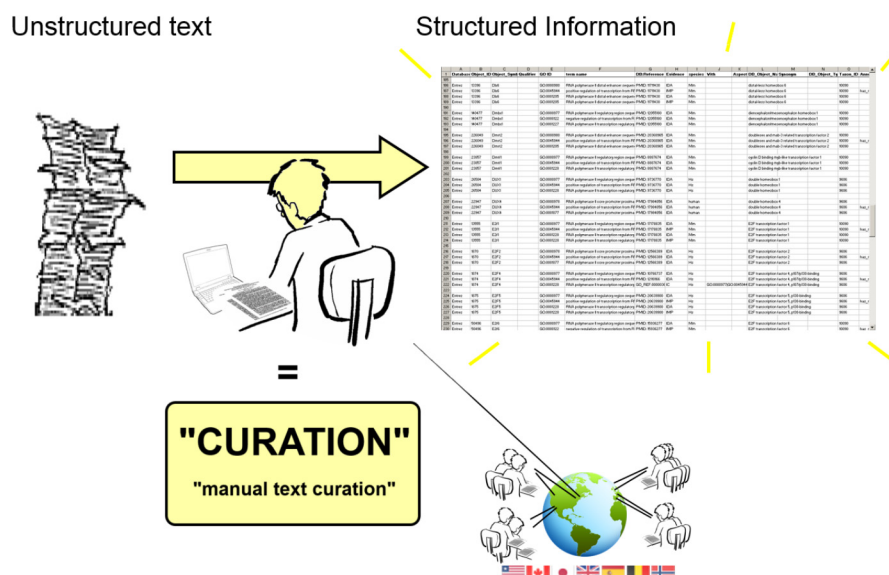
1. The current situation

Most biomedical research results culminate in a publication which only ends up getting buried into an ever growing pile of publications. The content of these papers is functionally invisible for both man and machine: invisible for researchers, because even if you work in a small research niche you should ideally read thousands of papers to know the existing research you need to build upon; and invisible for text-mining algorithms, because information in Life Science papers is too complex for these to comprehend. Only people can deal with detailed nuances of multi-molecular interactions, rich experimental contexts, heterogeneous facts, and often a need to infer what is actually meant.

This is especially important for new approaches using Systems Biology and Systems Medicine, as these fields require holistic and accurate overviews of many molecular interactions and cross-scale observations.

Most people who need some structured overview in their research niche have to create it themselves, and engage, without realizing it, in curation. While reading a selection of relevant papers, they write down a list of useful facts or collect them in a spreadsheet. During a decade of experience with various research labs, we have seen these local curation efforts happen over and over again: e.g. on plant cell cycle and biomass (whole department effort, handicapped by current technology), or gastric cancer network regulation (NTNU research group).

Most of these locally curated facts, however, get lost again and die at the end of a project.



Although institutional efforts by professional curator groups can efficiently carry out semi-automated expert curation strategies, they cannot keep up with the increasing amount of Life Science literature. Over a million new papers get published each year, and only a small percentage of selected information

gets structured. Most current curation projects, however, are in fact small-scale, as they are constrained by available *workforce* (limiting the curation volume) and by the supported *scope* (diversity) of content that spreadsheets can easily take in.

2. The key problem

There exists no easy and scalable tool to capture all of Life Science's extremely diverse and flexible (context-rich) information. Instead, a variety of curation approaches are used, often limited to specific types of facts, and with little possibility to capture details about context. Just look at two examples:

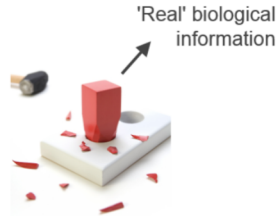
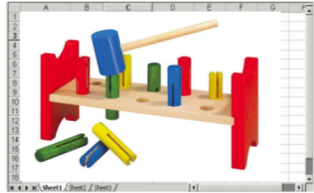
To capture details about a biological function with **spreadsheets**, one would not just want to describe e.g. "A activates B", but rather "A, *bound to D, transiently* activates a *modified form* of B, and only under certain *conditions*, as shown in *cell type X of species Y*, using *experiment Z*". For every type of essential context, you need an extra column. Such an approach soon becomes impractical as there is such heterogeneity of information in research results that a spreadsheet cannot capture it without adding an infinite number of columns. A curator would have to drop essential context, or abuse the form's semantics, or define a 'column X' that becomes a dump of free-text details that remains computationally opaque.

A more flexible format to capture diverse facts is a **controlled language**. However, consider "John eats chicken with fork". The language will have built-in rules that determine how sequences of words are interpreted. As "with" follows "chicken", the computer may capture this fact as "a chicken with a fork", instead of "eating with fork". Extrapolating this to biological sentences of many more words, it is easy to see the learning curve is extremely steep and only few people will take advantage of such languages' power.

The 'Curation Struggle'



1. Spreadsheets: are inflexible



➡ SHALLOW extraction

Most information
doesn't FIT
=> Few details curated

2. Formal Languages: are hard to learn

John eats Chicken with Fork



➡ FEW facts: experts only

HARD to explain to
computer
=> Few papers curated

3. Our new curation paradigm

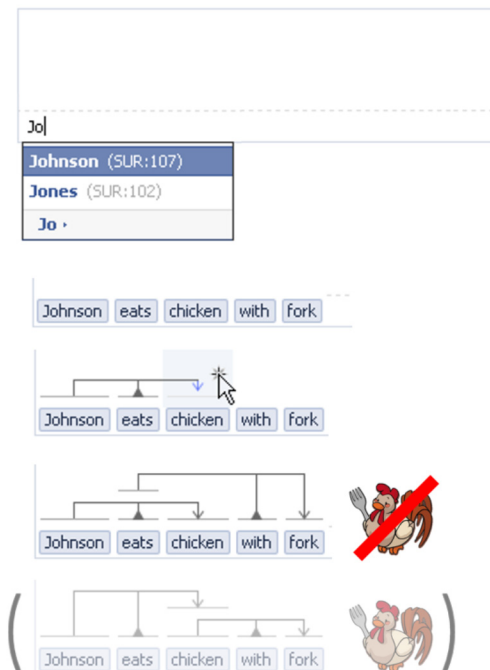
How can we enter any collection of words as an 'information statement' into software, and then make it clear how the words are connected, using only a minimal set of rules as building blocks, used over again?

For this, we developed *VSM: 'Visual Syntax Markup'*. Annotating with VSM is both 1) extremely flexible to capture diverse and rich information, solving the 'spreadsheet' problem, and 2) intuitive to understand and easy to use, solving the 'controlled language' problem.

A simple example is shown in the figure below. VSM uses identifier-linked *terms*. Every word or word-group, like "Open Science", must be linked to a semantic ID to make its meaning unambiguous. SciCura's prototype helps people achieve this via an autocomplete list with terms from ontologies and dictionaries, from a back-end database. The example uses a fictional dictionary. A curator chooses a name from this dictionary and adds words from lists of verbs, prepositions and nouns to construct a sentence with 5 'terms', each disambiguated. Next, the curator adds structure, simply by identifying *triples* in the sentence and marking them with a '*trident*' connector. The user-interface enables a curator to do this with three simple clicks: one above each of a set of terms in a triple, identifying the subject, the relation, and the object. The first trident specifies that an individual by the name of Johnson eats an animal of the species 'chicken'. The next trident clarifies that this "eating" is performed with a "fork".

The **key thing to understand about VSM** is the resulting semantics: every time a new trident is added, all connected terms receive extra context. For example, “eat” becomes “the eating of the chicken by Johnson”. Likewise “chicken” becomes “the chicken eaten by Johnson”, etc. With each concept becoming specific, one can refer to it, and attach further context like: “the eating of .. by ..”, that happens with fork. Like this, one can **keep adding further, nested context for every term**.

The example shows also that the tridents can be used to clarify the alternative interpretation: the imaginary case where a chicken holding a fork is being eaten. A demo video that illustrates this example step-by-step is attached as a media file to our submission.



We designed VSM so it can express just about everything a natural language can, by using three (and a half) connector types: the trident / dident; the list connector; and the co-reference connector that can also refer to terms in other sentences. See figure:

VSM: Quick tutorial

1. Trident



Stacking tridents:

= adding more *context* to each term



1b. Dident



2. List

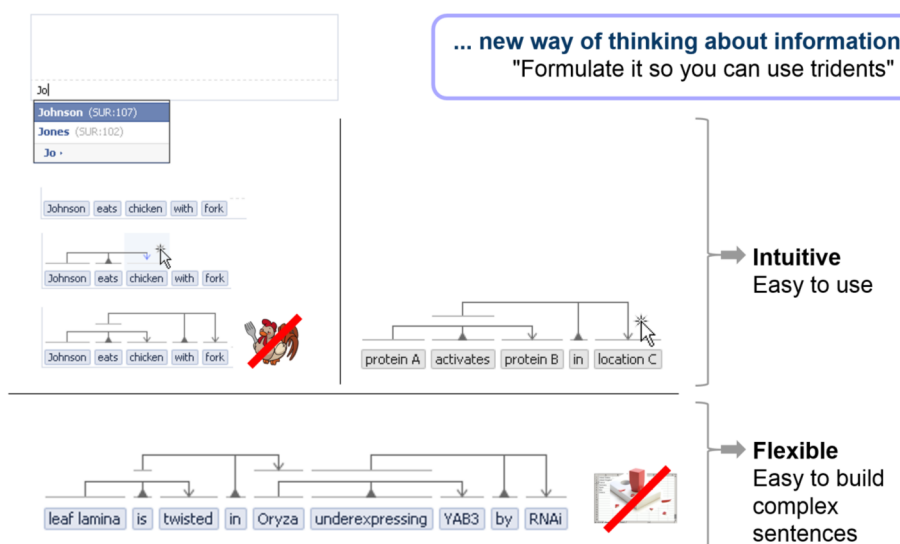


3. Co-reference



At the bottom of the next figure, an actual fact was built by repeated use of trident connectors adding more detail: “leaf lamina” (the blade of a leaf, an ID-linked term from Plant Ontology) was observed to be “twisted” (PATO term), and this being twisted “pertains to” (=ID-linked synonym of equivalent preposition “in”) the species “Oryza” (rice), underexpressing some gene, using RNA-interference (it can be crucial to know which experiment was used to assess a fact’s reliability). Such examples quickly become too complex for the spreadsheet method, whereas VSM can capture the biological context.

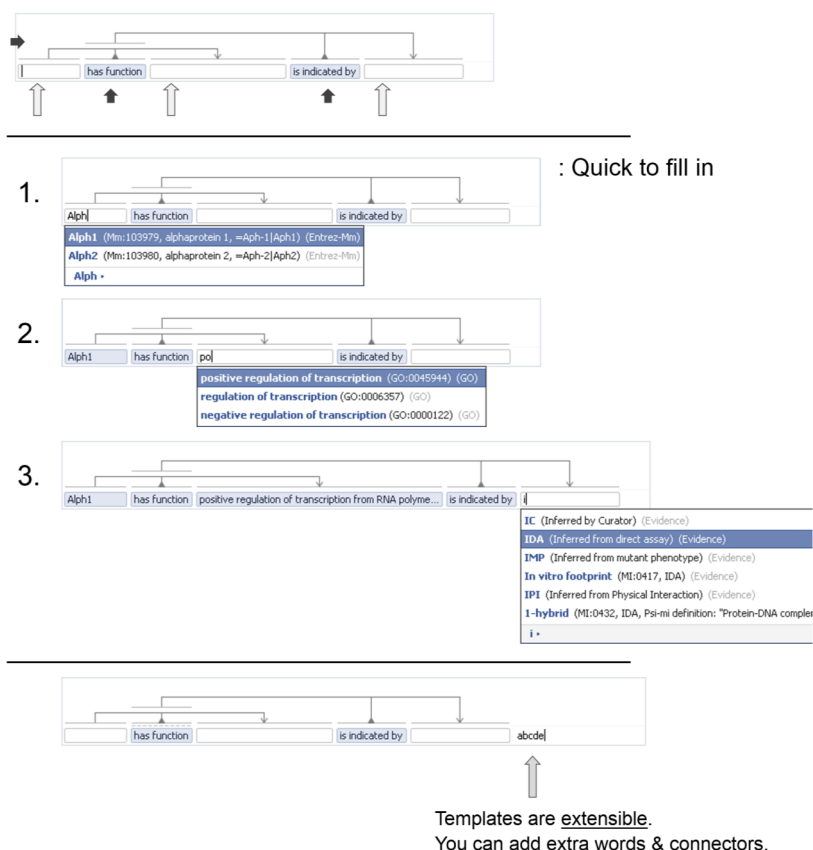
VSM = Flexible & easy method, powerful curation



To be clear: these sentences, or ‘information statements’, are *not* verbatim copied from a paper. They are reformulated by a user according to their interpretation of some experiment, and represent a disambiguated version of informational essence. In addition, tridents are *not* syntax parse trees generated by algorithms, but are manually added. Although in theory the results of text-mining may help to build such a sentence, this can be investigated later in the development of SciCura.

To simplify the entry of information statements using VSM we allow the use of *VSM-templates*, see figure below. These are VSM-statements with some pre-defined terms and connectors that are specific for the information type that should be recorded, with only some additional empty boxes to fill in. These can routinely be used for bulk curation. Yet unlike spreadsheets, empty boxes offer a selection of entity terms or ontology terms linked to their respective IDs. This means that no error-prone identifier-copy/paste steps are needed; only typing a few letters and accepting a suggested term by a mouse-click will do. In addition, templates are extensible: one can still add more words and link them up (avoiding the dreaded ‘column X’) to specify more context.

Or using **templates**:



4. SciCura-based community curation

With a flexible and simple curation method that enables a research community to curate whatever they need, we can develop SciCura: a web-based platform for community curation.

Our vision is that after several years it would have one page per paper, holding a VSM-sentence-based *digital abstract* of that paper's detailed findings. This would be similar to a Wikipedia-page, except that the content is fully semantic: the meaning of every term and complex sentence is clearly understood by software.

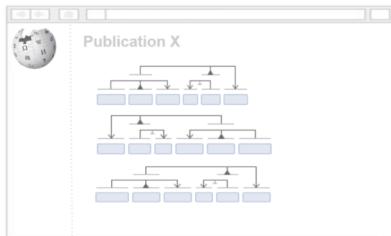
VSM enables SciCura web-platform

VSM method:  **Flexible** => you can curate **anything**
Flexible & easy curation method

Easy to use => **many** people can curate

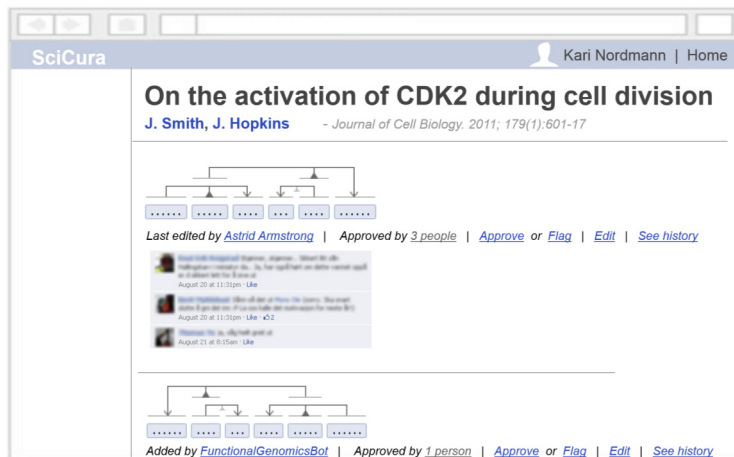


SciCura web-platform: **Large-scale curation**



Co-operation based web-platform
to create **Structured Summaries**
of Life-Science publications

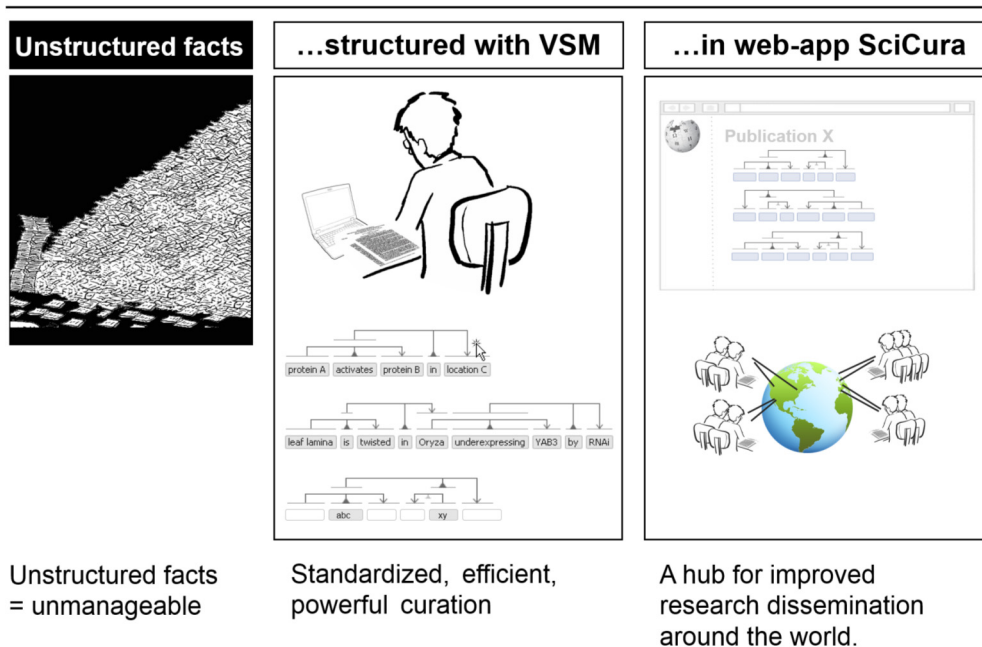
A fully social web platform would also support features like: commenting to discuss interpretations (a real curator need); an edit history; flagging of problems; and reviewing/approving of statements.



As a comparison: on Wikipedia, 90% just reads, 9% correct typos, and 1% writes content. Similarly, SciCura can grow when most people just fill templates, and some people use VSM's full power and design templates for others.

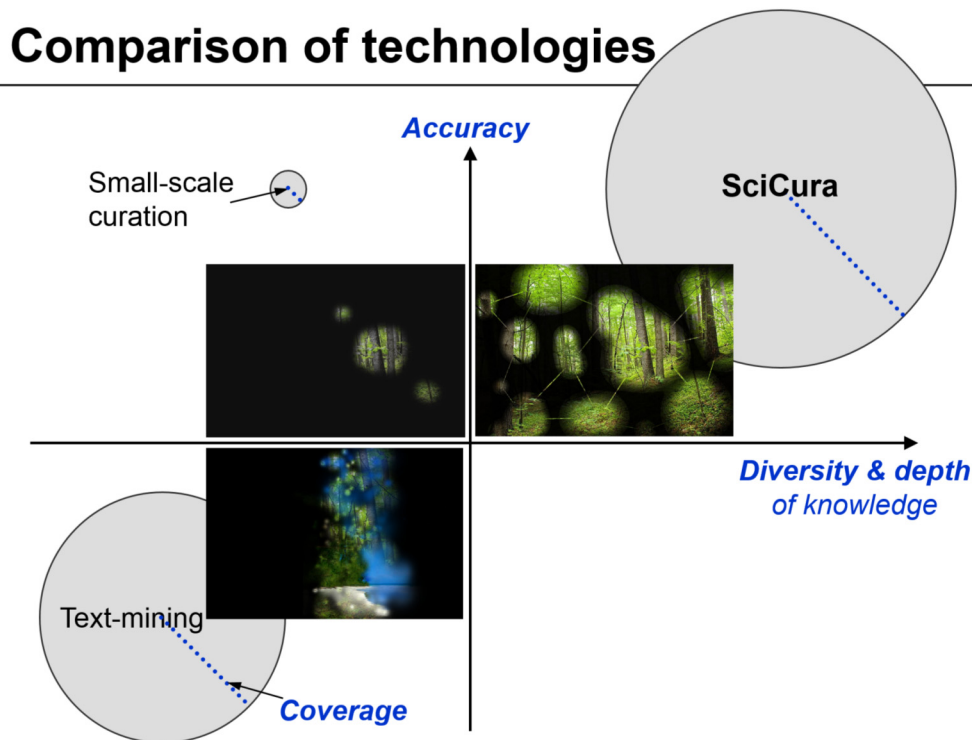
In summary: diverse, unstructured facts can be structured with VSM. This enables the web-app SciCura to open up previously opaque science results to the world, in a standardized and scalable way.

Summary so far



5. Positioning

A comparison to the current state-of-the-art is shown in the cartoon below. While small-scale curation is accurate and text-mining can boast coverage, both are still limited in scope. Meanwhile, SciCura's has potential for heterogeneous, crowd-sourced manual curation. Small-scale curation may be like walking in a dark forest with a narrow focused flashlight; text-mining is like using a magic lantern that shows many things that aren't really there (errors); SciCura could foster a global community curation effort of high quality and virtually limitless diversity and depth, increasingly showing the full details of the forest.



6. Impact

The impact of SciCura will increase with growing numbers of users. What incentives could there be to get SciCura adopted by a wide community of curators? In our earlier curation experiments, we directly integrated scientists' curations into a browsable overview (see our viewer OLSVis at ols.wordvis.com). The fact that curators could immediately see the results of their work in the context of what else had been entered by their colleagues proved to be an extremely useful 'carrot' for participants.

More generally, SciCura will give scientists access to the sum of curated knowledge through an API, and enable user-written output scripts, plugins, or even a shareable App repository. These additions may elevate SciCura to a community development platform much like the popular Cytoscape platform (cytoscape.org).

Exports of information to established databases (IntAct, the GO database, other) of the Knowledge Commons will further increase the impact. Several export possibilities are already in place (see section 7), and future export in the semantic web formats RDF and OWL will further the spread of information available through SciCura, including the use of graph searches. OWL exports will also facilitate the integration with the OWL editor Noctua (noctua.berkeleybop.org), and pave the way to apply DL reasoning to the collection of curated facts.

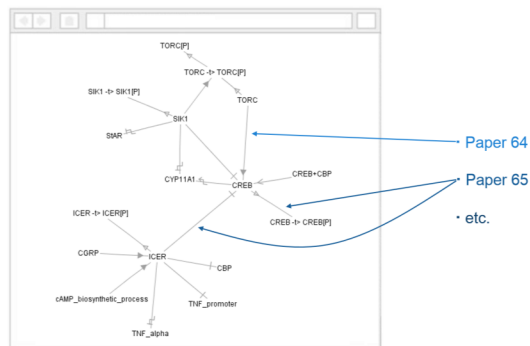
SciCura will also impact the development of text mining algorithms. A collection of structured facts linked to original papers or even sentences is an invaluable text corpus for training. Also, by cross-referencing each VSM statement to the original paper through PubMed IDs, platform users will be able to consult the original papers, which serves as an advertisement for journals and publishing houses.

Research also becomes more cost-efficient. Less duplicated effort of re-interpretation of complex results will happen, and relevant information becomes quicker to find.

Other benefits are shown in the figure below.

In summary, SciCura can be a **transformative project to open up knowledge**, making it computationally available and effectively accessible for Systems Biology, Systems Medicine and Precision Medicine. It can streamline access to open knowledge, deeply understood by intelligent machines. SciCura's growth depends on funding and publicity, making people contribute with code, ideas, and curated knowledge.

- Integrated Visualisation (=Direct benefit for curating)



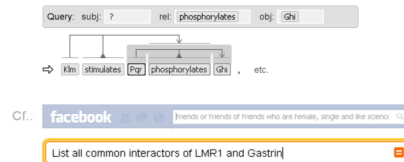
- Semantic search



- Information Management (sorting your information)



- Graph Search in the Life Sciences



- Systems Medicine / Agro / Biotech industry



- Text-mining development



7. A First Use Case and Prototype

The current prototype platform was developed in collaboration with volunteer curators at NTNU. Their focus is the curation of all DNA binding transcription factors (DbTFs) of human, mouse and rat; and a curation protocol was developed in collaboration with members of the Gene Ontology Consortium and curators of the Mouse Genome Database at the Jackson laboratories.

We use an agile, iterative strategy while building the prototype's functionality, in continuous discussion with curators. We prioritize the features that are needed first, while we maintain a balance to keep the web app *generic* enough for next use cases. The SciCura prototype currently features one *curation room* where six people can simultaneously curate DbTFs, their regulation effect, target gene, and experimental details.

We loaded the back-end database with various dictionaries and ontologies needed for the task. As these can be incomplete (e.g. Cell Type), new terms can be added ad-hoc and then reused, and are flagged.

Each curation statement is linked to the original paper, here simply via PubMed ID. We added a place to put the 'original sentence' from which information was extracted, even though facts were often inferred from several parts of text and figures. Curators started using this place to communicate: comments on

interpretation, flagging if re-checked. This shows a user need to develop social and other features around curated sentences.

1) An early version of the curation room:

Selection of templates
(short, early-stage templates)

Templates designed with domain experts

Multi-user, cooperative

2) One of five current templates, larger to hold additional details; and an autocomplete list:

HeLa cell (BTO:0000567, Human cervix carcinoma established from the epitheloid cervix carcinoma of a 31-year-old

HeLa GFP-histone H2B cell (BTO:0003602, HeLa cells expressing a GFP-histone H2B protein.) (BTO)

HeLa-80 cell (BTO:0003180, A strain of HeLa cells that proliferates under 80% O2, termed HeLa-80, has been deriv

HeLa-229 cell (BTO:0005030, Human epitheloid cervix carcinoma cell line from a black female.) (BTO)

HeLa-LTRHIV-1-Luc cell (BTO:0005038) (BTO)

HeLa-MAGI cell (BTO:0003588, CD4 positive HeLa cell line that contains an integrated HIV-1 promoter.) (BTO)

HeLa-MAGI-CCR5 cell (BTO:0003587, A HeLa-CD4 cell line that expresses CCR5 and that has an integrated copy c

HeLa-S3 cell (BTO:0000568, Human cervix carcinoma cell line is a subclone of its parent HELA derived in 1955.) (B

chela (BTO:0002772, A pincerlike claw of a crustacean or arachnid, such as a lobster, crab, or scorpion.) (BTO)

hela -

3) A selection from 5000+ curated sentences; under an input box with 5 selectable templates:

The screenshot displays the scicura.org web interface. At the top, there is a search bar with the URL 'scicura.org' and various filters like 'Show all text', 'in all period', 'and all template', 'by anyone', 'as a', 'curator', 'Latest', '15 after', '15', 'Apply', 'Reset', '« Prev', 'Next »', 'Output', 'Steven Verduyse', and 'Log Out'. Below the search bar, there is a sidebar with a 'No template' button and a list of templates: 'TF-TG - MF', 'TF-TG - BP', 'TF2-TG - MF', 'TF2-TG - BP', and 'TF-TF'. The main content area shows five example sentences, each with a corresponding trident structure. The sentences are: 1) 'MYOD1 is represented by recombinant TF with His tag ... has function RNA polymerase II regulatory region sequence-spe... ACACB is represented by TFBS, authentic, competi... is indicated by electrophoretic mobility shift assay reported in Fig 8A in in vitro'. 2) 'NAC1 is represented by recombinant TF in vitro trans... has function RNA polymerase II regulatory region sequence-spe... gene is represented by TFBS, consensus is indicated by electrophoretic mobility shift assay reported in Fig 3 in in vitro'. 3) 'NFX1 is represented by recombinant TF without tag ... has function RNA polymerase II regulatory region sequence-spe... HLA-DRA is represented by TFBS, authentic, competi... is indicated by electrophoretic mobility shift assay reported in Fig 4A in in vitro'. 4) 'Onec43 is represented by overexpression of recomb... is involved in positive regulation of transcription from RNA polym... Foxa2 is represented by TFBS, authentic, promoter is indicated by reporter gene assay reported in Fig 3 in HEK-293 cell'. 5) 'Onec43 is represented by nuclear extract from cells or ... has function RNA polymerase II core promoter proximal region s... Foxa2 is represented by TFBS, authentic, promoter,... is indicated by electrophoretic mobility shift assay reported in Fig 2B in COS-7 cell'. Each example sentence is followed by a trident structure diagram and a PubMed ID.

Our users testify they don't find the trident structures difficult at all. In fact, they are excited about all the complexity they can now handle in such an elegant manner, allowing them to focus on the biology without worries about error-prone entry sheets.

The current prototype can generate output of relevant parts of information to established databases like GO (geneontology.org) and IntAct (www.ebi.ac.uk/intact). We built a graph-query API that can fish out any matching pattern or sub-pattern of terms and connectors in all VSM-sentences. Then we wrote output-scripts to fish out, combine and reformat relevant sets of information.

scicure

scicure

Select output type:

Combined GO&IntAct	GO (GPAD)	IntAct (MITAB)	Plain
Combined GO&IntAct + save to file	GO (GPAD) + save to file	IntAct (MITAB) + save to file	Plain + save to file

Calculating InAct (MITAB) output (for MF):										Calculating GO (GPAD) output (for MF & BP):																																																																																																																																																																																																																																																																																																																																																																																																											
> Querying for template 1: 1270 matches. > Reads:										> Querying for template 1: 1270 matches. > Reads:																																																																																																																																																																																																																																																																																																																																																																																																											
id	id1	id2	id3	id4	id5	id6	id7	id8	id9	id10	id11	id12	id13	id14	id15	id16	id17	id18	id19	id20	id21	id22	id23	id24	id25	id26	id27	id28	id29	id30	id31	id32	id33	id34	id35	id36	id37	id38	id39	id40	id41	id42	id43	id44	id45	id46	id47	id48	id49	id50	id51	id52	id53	id54	id55	id56	id57	id58	id59	id60	id61	id62	id63	id64	id65	id66	id67	id68	id69	id70	id71	id72	id73	id74	id75	id76	id77	id78	id79	id80	id81	id82	id83	id84	id85	id86	id87	id88	id89	id90	id91	id92	id93	id94	id95	id96	id97	id98	id99	id100	id101	id102	id103	id104	id105	id106	id107	id108	id109	id110	id111	id112	id113	id114	id115	id116	id117	id118	id119	id120	id121	id122	id123	id124	id125	id126	id127	id128	id129	id130	id131	id132	id133	id134	id135	id136	id137	id138	id139	id140	id141	id142	id143	id144	id145	id146	id147	id148	id149	id150	id151	id152	id153	id154	id155	id156	id157	id158	id159	id160	id161	id162	id163	id164	id165	id166	id167	id168	id169	id170	id171	id172	id173	id174	id175	id176	id177	id178	id179	id180	id181	id182	id183	id184	id185	id186	id187	id188	id189	id190	id191	id192	id193	id194	id195	id196	id197	id198	id199	id200	id201	id202	id203	id204	id205	id206	id207	id208	id209	id210	id211	id212	id213	id214	id215	id216	id217	id218	id219	id220	id221	id222	id223	id224	id225	id226	id227	id228	id229	id230	id231	id232	id233	id234	id235	id236	id237	id238	id239	id240	id241	id242	id243	id244	id245	id246	id247	id248	id249	id250	id251	id252	id253	id254	id255	id256	id257	id258	id259	id260	id261	id262	id263	id264	id265	id266	id267	id268	id269	id270	id271	id272	id273	id274	id275	id276	id277	id278	id279	id280	id281	id282	id283	id284	id285	id286	id287	id288	id289	id290	id291	id292	id293	id294	id295	id296	id297	id298	id299	id300	id301	id302	id303	id304	id305	id306	id307	id308	id309	id310	id311	id312	id313	id314	id315	id316	id317	id318	id319	id320	id321	id322	id323	id324	id325	id326	id327	id328	id329	id330	id331	id332	id333	id334	id335	id336	id337	id338	id339	id340	id341	id342	id343	id344	id345	id346	id347	id348	id349	id350	id351	id352	id353	id354	id355	id356	id357	id358	id359	id360	id361	id362	id363	id364	id365	id366	id367	id368	id369	id370	id371	id372	id373	id374	id375	id376	id377	id378	id379	id380	id381	id382	id383	id384	id385	id386	id387	id388	id389	id390	id391	id392	id393	id394	id395	id396	id397	id398	id399	id400	id401	id402	id403	id404	

The next development phase is adding the following elements:

We enable two use cases:

15

A2) Deep curation: a pilot in mouse information curation (Mouse Genome Informatics, JAX). Here the challenge is to develop a flexible template building utility, allowing curators to build their own templates for genetic, genomic and phenotype curation.

B) Integration of ontologies, biological dictionaries, term suggestion service TermGenie (Mungall, BerkeleyBOP):

We will integrate the necessary dictionaries from NCBI Gene, UniProt and NCBI BioPortal into the system. As there may often be a need for new terms/concepts, we will allow users to define candidate terms with a minimal definition and basic categorization, after which candidate terms will be processed/authorised by the TermGenie application (go.termgenie.org).

C) SciCura exports (in collaboration with Sandra Orchard, IntAct, EMBL-EBI; Mungall, Noctua, BerkeleyBOP):

We will further develop the export of facts to the IntAct database. We will in addition enable export to Noctua (OWL format), allowing further refinement of knowledge statements and DL reasoning to classify the information and produce new assertions.

Our work is currently funded by three grants, guaranteeing that by the end of 2016 we will have a platform that can be opened to the general public and open-sourced on GitHub. During this year, our cooperations will help us prioritize between many features we could build, such as:

- multiple curation rooms;
- prefill with statements from text-mining (or from an input-script);
- the infrastructure to support user-defined plugins that query the SciCura content (e.g. to enable custom visualisations). This will set the stage for further community development of the system;
- public/private curations rooms for curation-before-publication, similar to GitHub;
- advanced search based on reasoning.

For continued outreach, we will educate people in use of the platform; present it at conferences, (social) media, and to key biomedical scientists; publish written and video tutorials, and offer training sessions at conferences like ISMB and ECCB, and the International Society for Biocuration.