

Open-Source Data Science Projects at eLife

Daniel Ecer
22th November 2019



eLife

What is eLife?

hhmi | Howard Hughes
Medical Institute



*Knut and Alice
Wallenberg
Foundation*

- A non-profit backed by **research funders** to drive reform in research communication
- We invest heavily in open-source technology development and innovation **on behalf of the community**



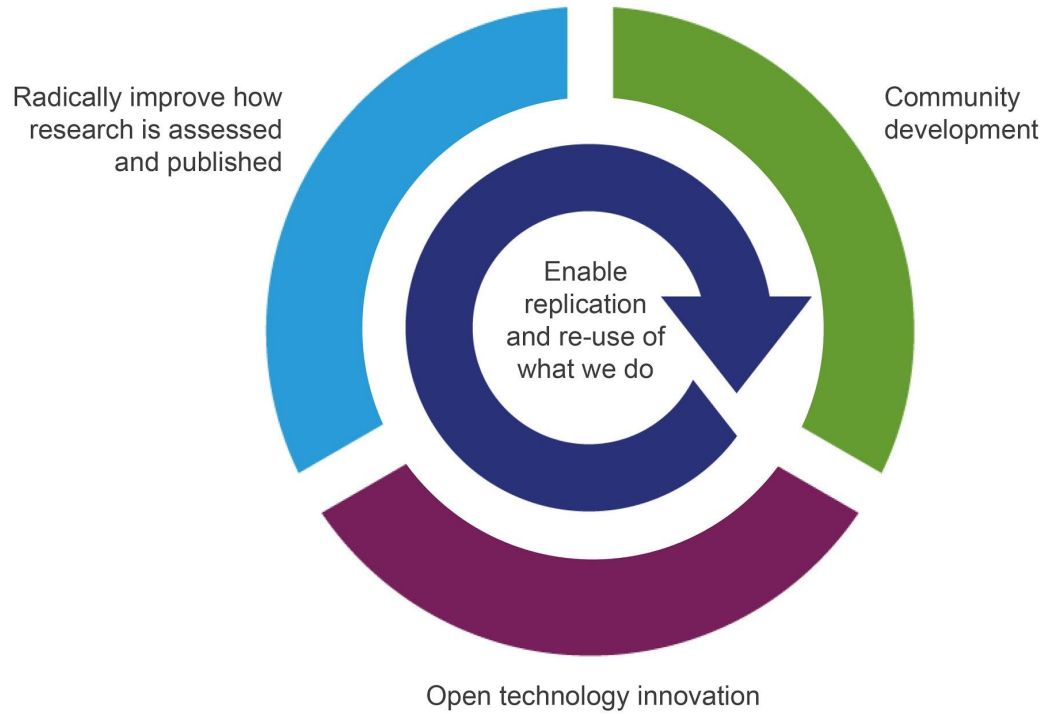
Helping scientists **accelerate discovery** by
operating a platform for research **communication**
that encourages and recognises **the most**
responsible behaviours in science.

e L I

eLife's motivations

- Leverage the **power of web technology** to accelerate research and discovery
- **Support open-access** publishing
- Build a **community-owned infrastructure** for research communication





Data Science projects at eLife



Amy
Author

- Amy submits a manuscript to a journal, maybe not the first attempt
- Amy would like to reduce form filling
- **ScienceBeam** helps pre-populating fields
- **ScienceBeam** may also help to add a semantic structure early on, to make the manuscript more “accessible”



Hubert
Reviewing Editor / Editorial Staff

- Hubert has been assigned, as the reviewing editor
- Hubert now needs to assign reviewers
- **PeerScout** helps to reduce bias and extend pool to find peers at different stages (editors, reviewers etc)



Philip
Reader

- Philip wants to stay on top of his field of interest
- There are many published manuscripts
- With **PeerTax**, Philip can better incorporate the views of reviewers, when reading manuscripts
- While **Citation Sentiment** provides a better picture of why manuscripts were cited

ScienceBeam

Can you guess...?

1?

2?

3?

4?

5?

6?

CV model qualitative results

Input
(PDF page)



input

Journal of Experimental & Clinical Assisted Reproduction

Comparison of selected cryoprotective agents to stabilize zootic spermatids of human oocytes during cooling

Abstract

Background: This study compared the protective effect of selected cryoprotective agents (CPA) on the survival of human oocytes during cooling.

Methods: Human oocytes were collected from patients undergoing IVF treatment and subjected to cryopreservation in 100 oocytes with multiple oocytes at 0°C, 10°C and 20°C and subjected to 20°C to further stabilize oocytes. The oocytes that survived at 0°C, 10°C and 20°C after being cryoprotected with 1,2,3,4,5,6,7,8,9,10,11,12,13,14,15,16,17,18,19,20,21,22,23,24,25,26,27,28,29,30,31,32,33,34,35,36,37,38,39,40,41,42,43,44,45,46,47,48,49,50,51,52,53,54,55,56,57,58,59,60,61,62,63,64,65,66,67,68,69,70,71,72,73,74,75,76,77,78,79,80,81,82,83,84,85,86,87,88,89,90,91,92,93,94,95,96,97,98,99,100,101,102,103,104,105,106,107,108,109,110,111,112,113,114,115,116,117,118,119,120,121,122,123,124,125,126,127,128,129,130,131,132,133,134,135,136,137,138,139,140,141,142,143,144,145,146,147,148,149,150,151,152,153,154,155,156,157,158,159,160,161,162,163,164,165,166,167,168,169,170,171,172,173,174,175,176,177,178,179,180,181,182,183,184,185,186,187,188,189,190,191,192,193,194,195,196,197,198,199,200,201,202,203,204,205,206,207,208,209,210,211,212,213,214,215,216,217,218,219,220,221,222,223,224,225,226,227,228,229,230,231,232,233,234,235,236,237,238,239,240,241,242,243,244,245,246,247,248,249,250,251,252,253,254,255,256,257,258,259,260,261,262,263,264,265,266,267,268,269,270,271,272,273,274,275,276,277,278,279,280,281,282,283,284,285,286,287,288,289,290,291,292,293,294,295,296,297,298,299,300,301,302,303,304,305,306,307,308,309,310,311,312,313,314,315,316,317,318,319,320,321,322,323,324,325,326,327,328,329,330,331,332,333,334,335,336,337,338,339,340,341,342,343,344,345,346,347,348,349,350,351,352,353,354,355,356,357,358,359,360,361,362,363,364,365,366,367,368,369,370,371,372,373,374,375,376,377,378,379,380,381,382,383,384,385,386,387,388,389,390,391,392,393,394,395,396,397,398,399,400,401,402,403,404,405,406,407,408,409,410,411,412,413,414,415,416,417,418,419,420,421,422,423,424,425,426,427,428,429,430,431,432,433,434,435,436,437,438,439,440,441,442,443,444,445,446,447,448,449,450,451,452,453,454,455,456,457,458,459,460,461,462,463,464,465,466,467,468,469,470,471,472,473,474,475,476,477,478,479,480,481,482,483,484,485,486,487,488,489,490,491,492,493,494,495,496,497,498,499,500,501,502,503,504,505,506,507,508,509,510,511,512,513,514,515,516,517,518,519,520,521,522,523,524,525,526,527,528,529,530,531,532,533,534,535,536,537,538,539,540,541,542,543,544,545,546,547,548,549,550,551,552,553,554,555,556,557,558,559,560,561,562,563,564,565,566,567,568,569,570,571,572,573,574,575,576,577,578,579,580,581,582,583,584,585,586,587,588,589,590,591,592,593,594,595,596,597,598,599,600,601,602,603,604,605,606,607,608,609,610,611,612,613,614,615,616,617,618,619,620,621,622,623,624,625,626,627,628,629,630,631,632,633,634,635,636,637,638,639,640,641,642,643,644,645,646,647,648,649,650,651,652,653,654,655,656,657,658,659,660,661,662,663,664,665,666,667,668,669,670,671,672,673,674,675,676,677,678,679,680,681,682,683,684,685,686,687,688,689,690,691,692,693,694,695,696,697,698,699,700,701,702,703,704,705,706,707,708,709,710,711,712,713,714,715,716,717,718,719,720,721,722,723,724,725,726,727,728,729,730,731,732,733,734,735,736,737,738,739,740,741,742,743,744,745,746,747,748,749,750,751,752,753,754,755,756,757,758,759,760,761,762,763,764,765,766,767,768,769,770,771,772,773,774,775,776,777,778,779,780,781,782,783,784,785,786,787,788,789,790,791,792,793,794,795,796,797,798,799,800,801,802,803,804,805,806,807,808,809,810,811,812,813,814,815,816,817,818,819,820,821,822,823,824,825,826,827,828,829,830,831,832,833,834,835,836,837,838,839,840,841,842,843,844,845,846,847,848,849,850,851,852,853,854,855,856,857,858,859,860,861,862,863,864,865,866,867,868,869,870,871,872,873,874,875,876,877,878,879,880,881,882,883,884,885,886,887,888,889,890,891,892,893,894,895,896,897,898,899,900,901,902,903,904,905,906,907,908,909,910,911,912,913,914,915,916,917,918,919,920,921,922,923,924,925,926,927,928,929,930,931,932,933,934,935,936,937,938,939,940,941,942,943,944,945,946,947,948,949,950,951,952,953,954,955,956,957,958,959,960,961,962,963,964,965,966,967,968,969,970,971,972,973,974,975,976,977,978,979,980,981,982,983,984,985,986,987,988,989,990,991,992,993,994,995,996,997,998,999,1000.

target Target (expected) prediction

target

prediction

Prediction (actual)

input

Genetic Approach for the Fast Discovery of Phenazine Producing Bacteria

Abstract

Background: Phenazine is a class of natural products with diverse biological activities. The discovery of new phenazine-producing bacteria is a challenge for researchers. In this study, we developed a genetic approach for the fast discovery of phenazine-producing bacteria.

Methods: We used a genetic approach to identify phenazine-producing bacteria. The approach involved the use of a phenazine biosynthetic gene cluster as a probe to identify phenazine-producing bacteria. The results showed that the genetic approach was effective in identifying phenazine-producing bacteria.

Conclusion: The genetic approach is a fast and efficient method for the discovery of phenazine-producing bacteria. This approach can be used to identify new phenazine-producing bacteria and to study the biosynthesis of phenazine.

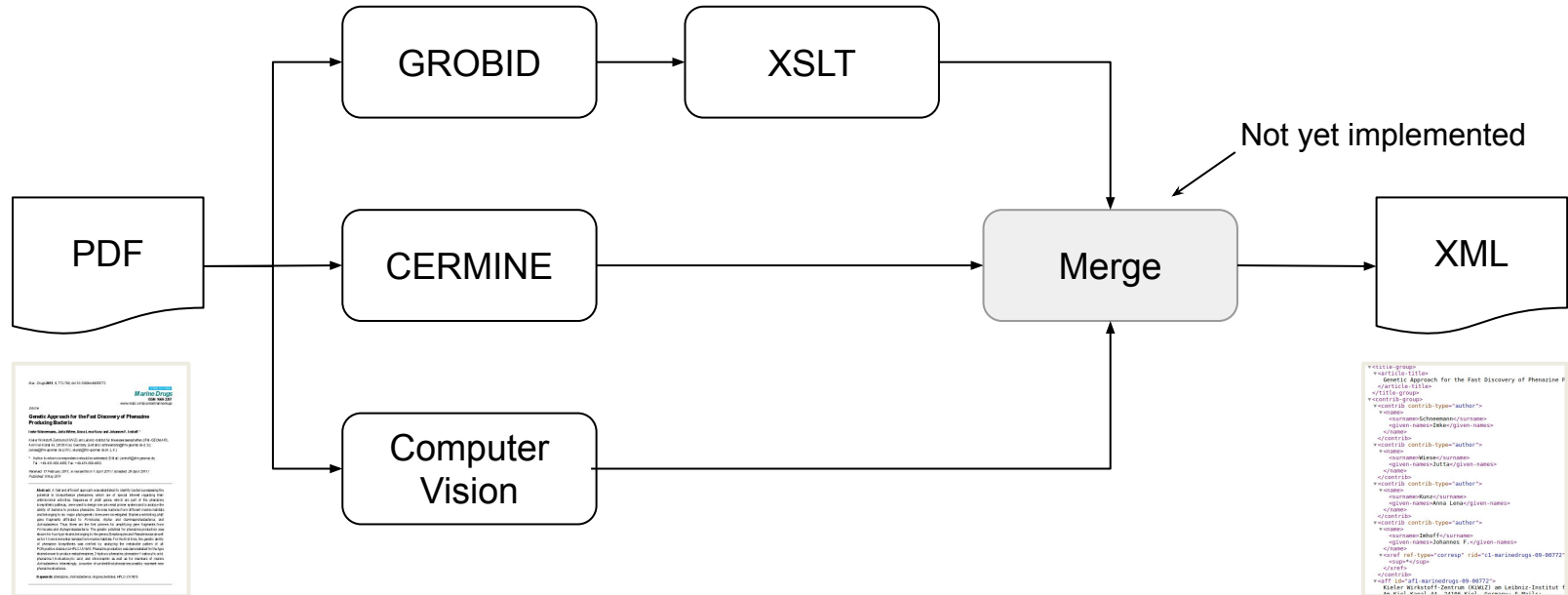
target

target

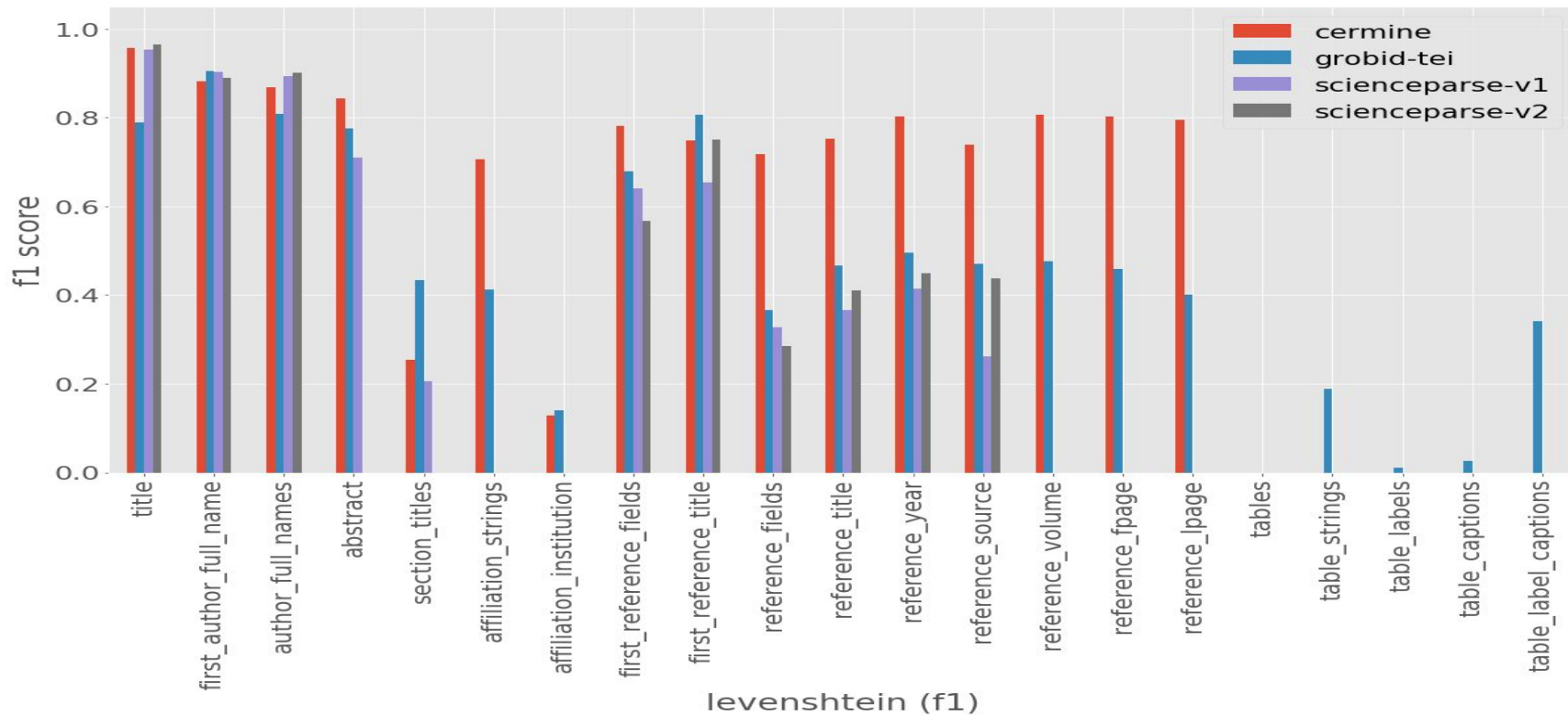
prediction

prediction

ScienceBeam - potential pipeline

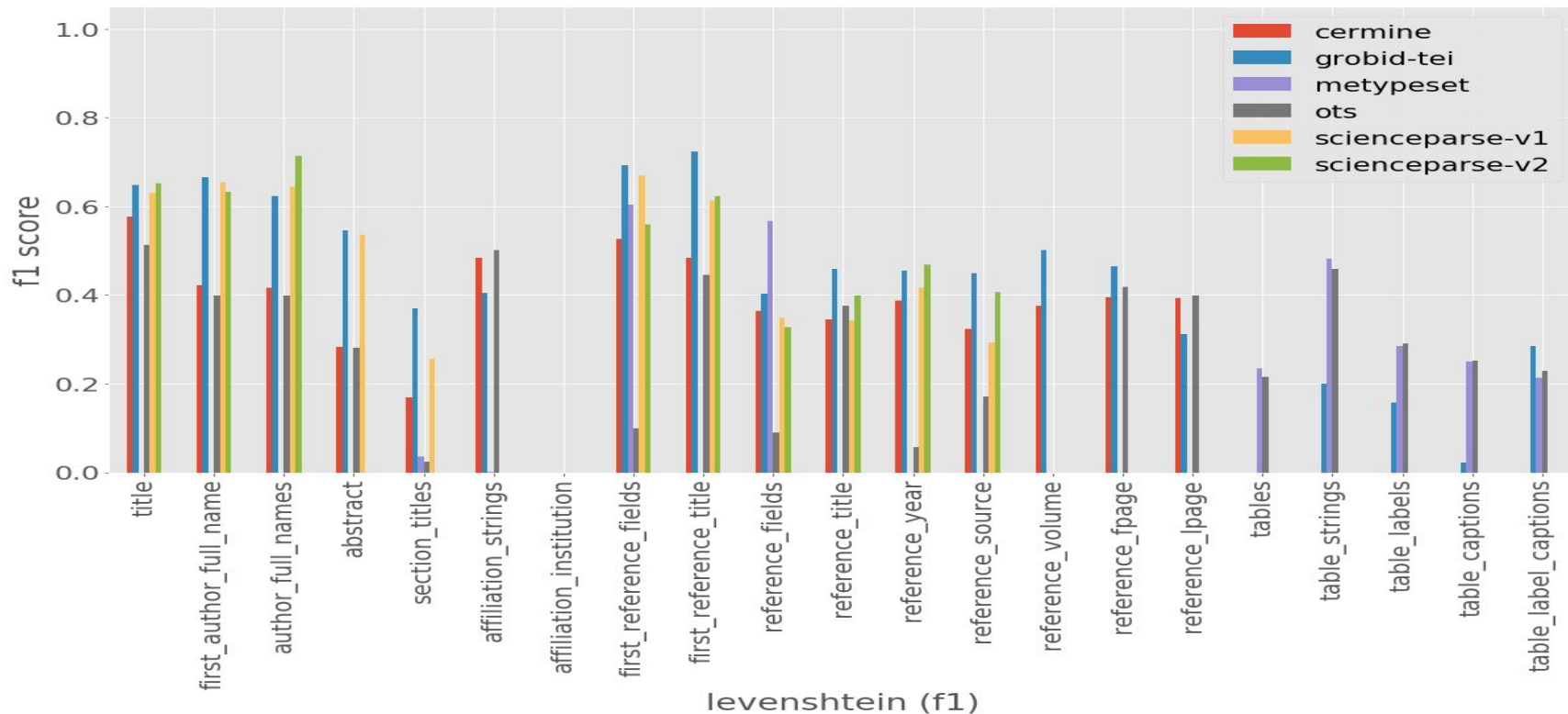


Evaluation - PMC_sample_1943

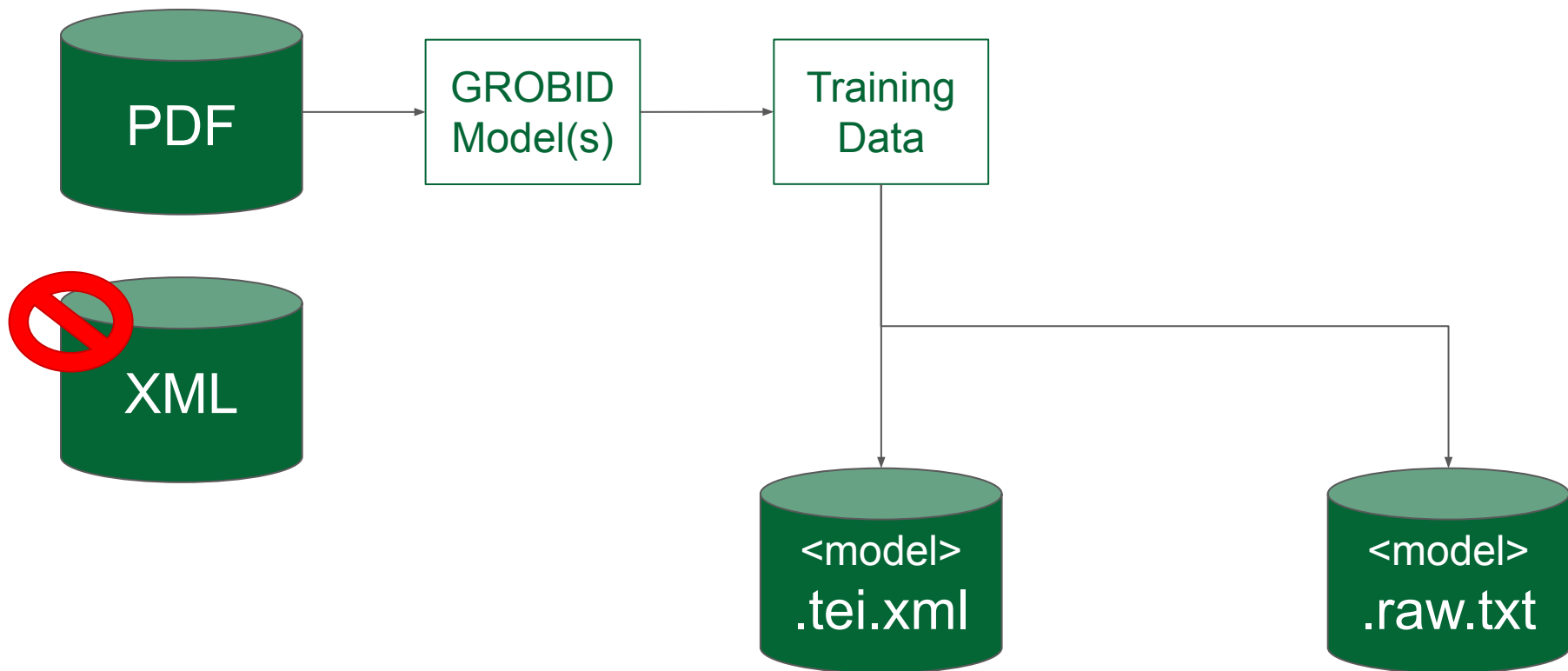


*CERMINE and ScienceParse are trained on PMC manuscripts

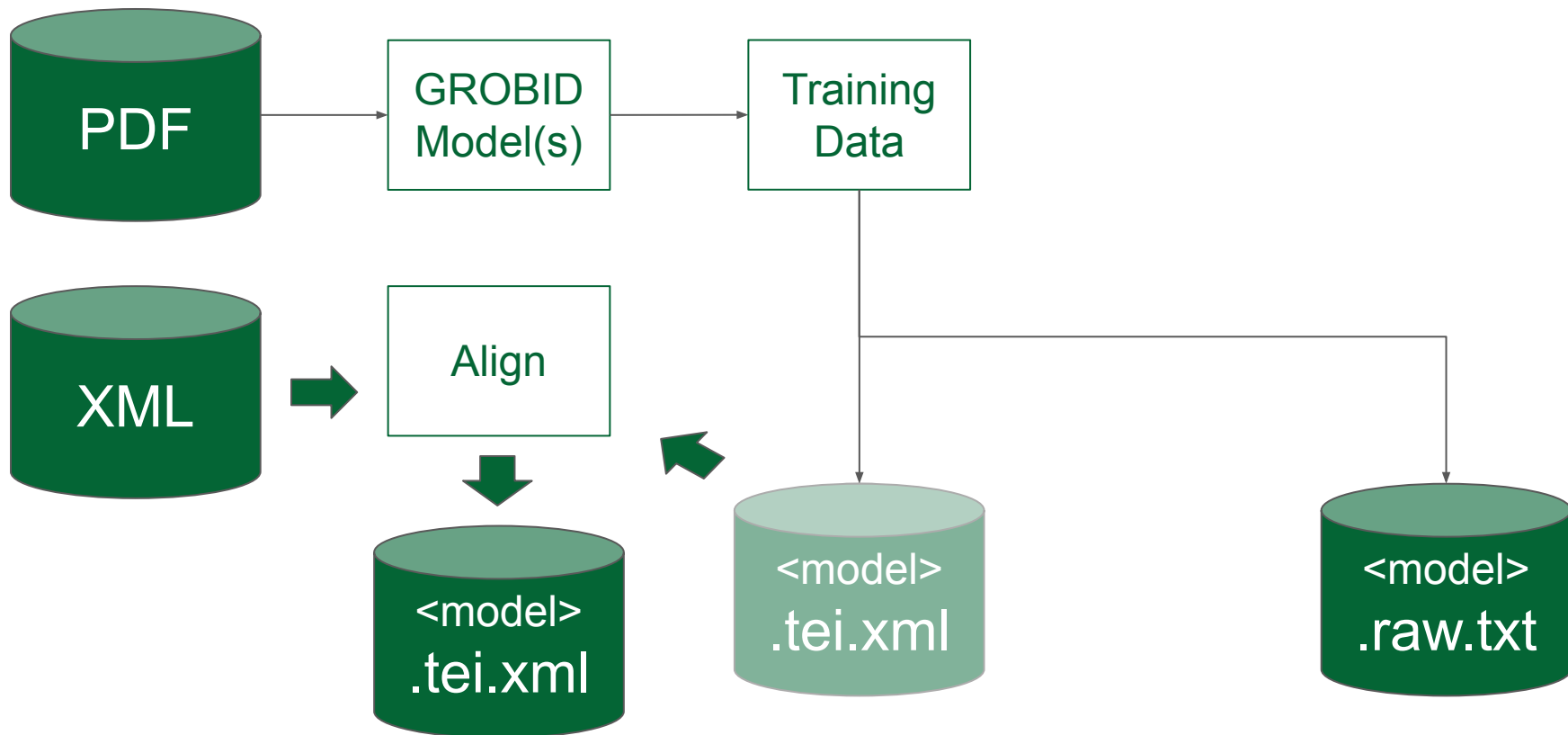
Evaluation - PKP coaction (Word)



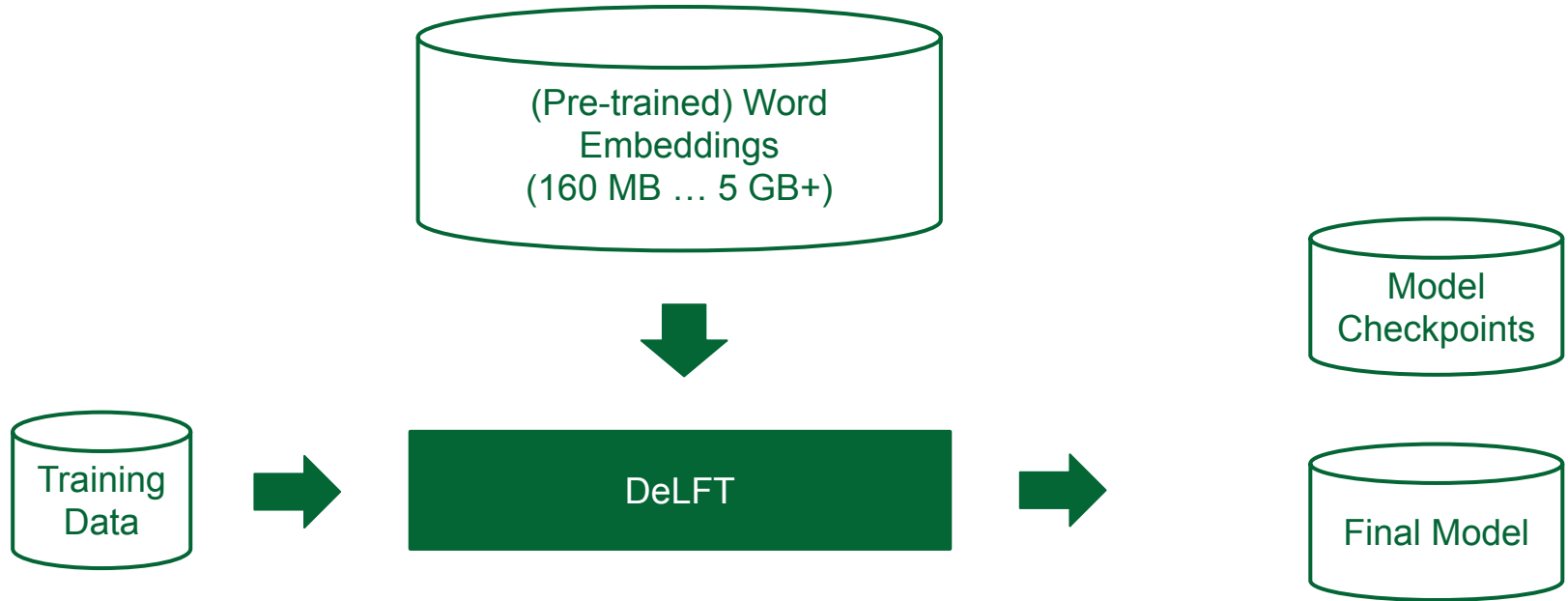
General GROBID Training Data Generation



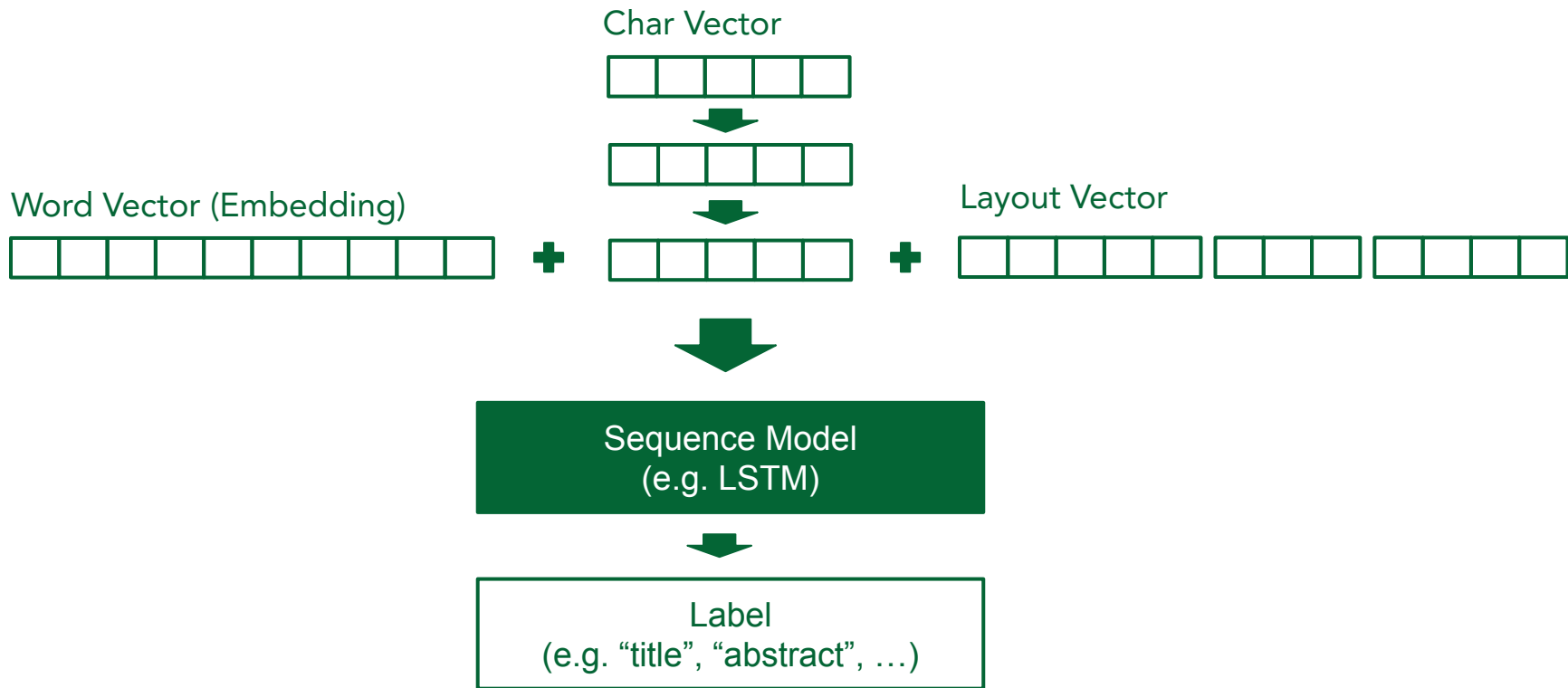
Use (JATS) XML to auto-annotate TEI training XML



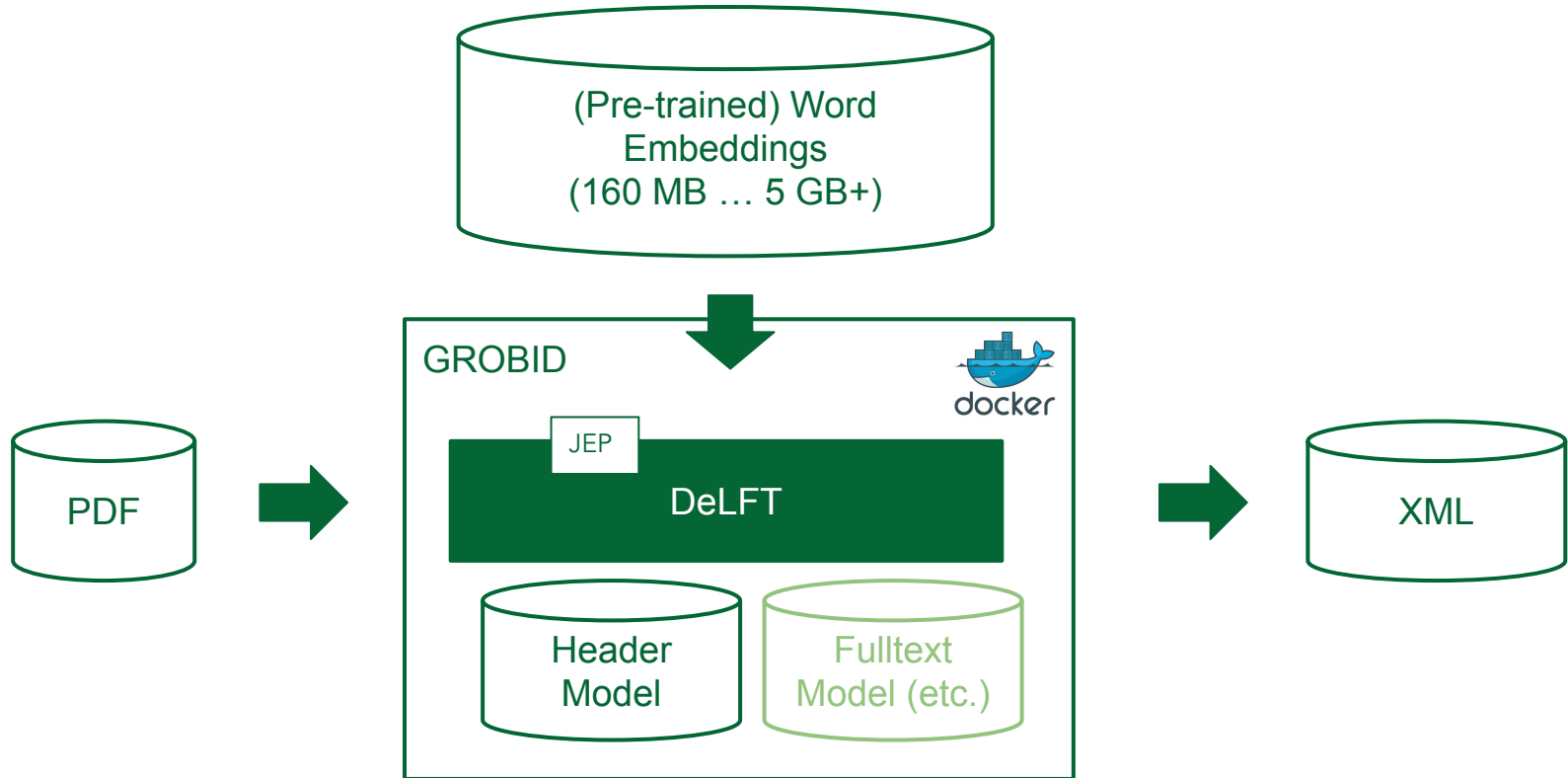
Training DeLFT Model



Adding layout features

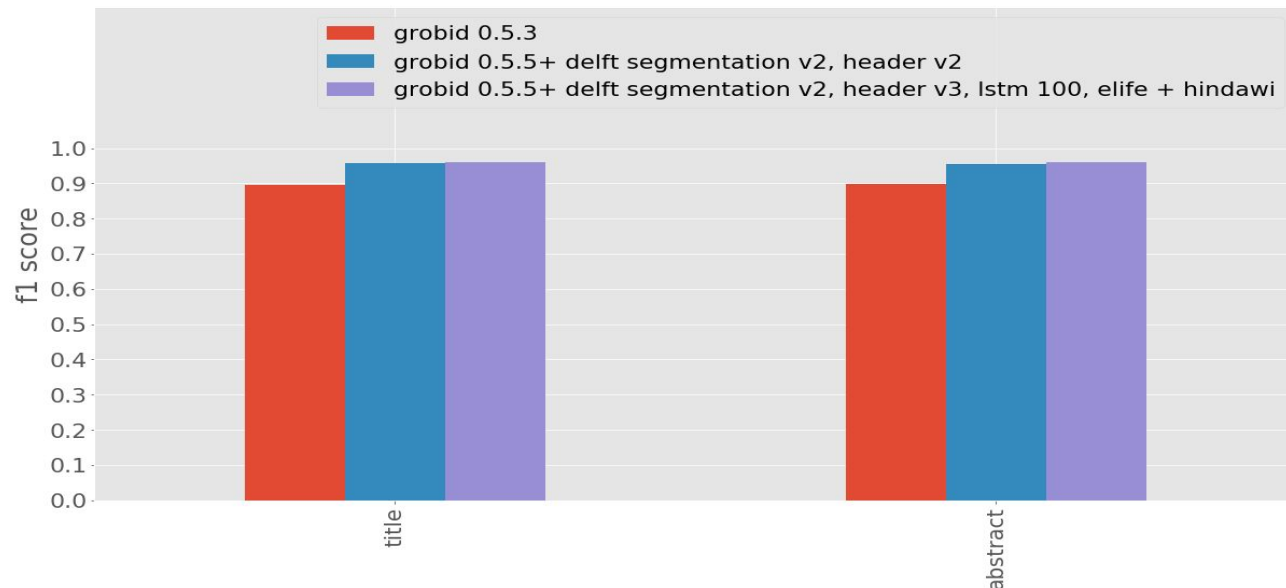


GROBID with DeLFT (in-progress)



Selected models, /wo header/footer, with abstract only

	grobid 0.5.3	grobid 0.5.5+ delft segmentation v2, header v2	grobid 0.5.5+ delft segmentation v2, header v3, lstm 100, elife + hindawi
title	0.895225	0.957947	0.960798
abstract	0.899202	0.954903	0.961373



eLife author submitted (PDF v2 /wo header/footer /w abstracts) - levenshtein (f1)

DAGs

Search:

		DAG	Schedule	Owner	Recent Tasks	Last Run	DAG Runs	Links
	<input type="checkbox"/> On	sciencebeam_convert	None	Airflow	8 1 39	2019-03-26 17:01	1 7	
	<input type="checkbox"/> On	sciencebeam_evaluate	None	Airflow	3	2019-03-26 17:16	1	
	<input type="checkbox"/> On	sciencebeam_evaluation_results_to_bq	None	Airflow	3 1	2019-03-26 17:20	1	
	<input type="checkbox"/> On	sciencebeam_watch_experiments	None	Airflow	5 1 4	2019-03-26 17:02	1 1	

Showing 1 to 4 of 4 entries

[Show Paused DAGs](#)

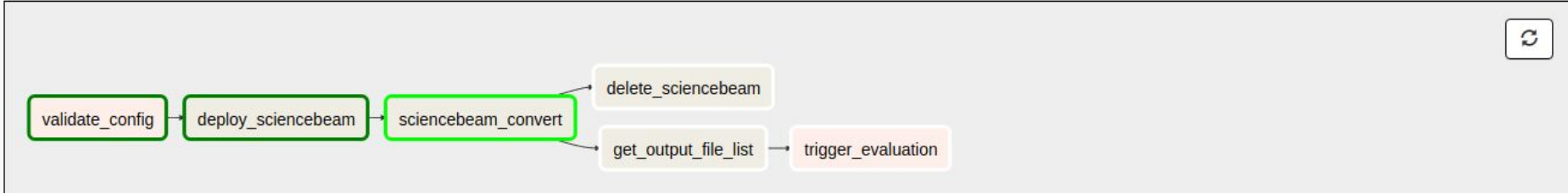
Graph View Tree View Task Duration Task Tries Landing Times Gantt Details Code Refresh Delete

running Base date: 2019-03-26 17:02:00 Number of runs: 25 Run: trig__2019-03-26T17:01:53.453502+00:00_elife-author-submitted-pdf-validation_grobid-tei-0.5.2 Layout: Left->Right Go

Search for...

BashOperator PythonOperator TriggerDagRunOperator

success running failed skipped rescheduled retry queued no status



¿y como generamos XML?

Semantic Extraction Working Group



eLIFE



PKP

PUBLIC
KNOWLEDGE
PROJECT

éruudit

+ algunos otros

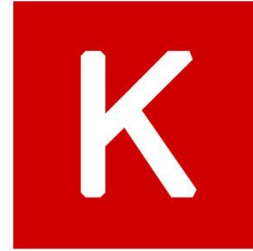
JUAN PABLO
ALPERIN

JUAN PABLO ALPERIN

@juancommander @pkp #scholcomm|lab



Some of the technologies used..

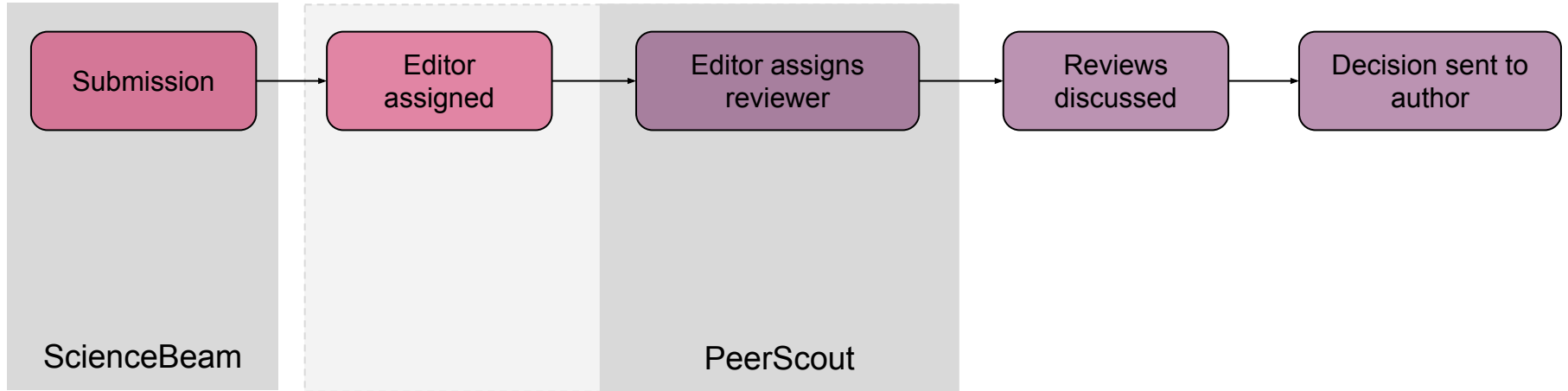


Conclusion

- Lessons learned:
 - Active community more important than higher initial score
 - Setting up training pipeline slow, compromise using Jupyter + Kubernetes for training and Airflow + Kubernetes for evaluation
 - Significantly improved performance using GROBID DL + layout features
- Next:
 - Train on larger data source
 - Improve other elements
 - Extend use-cases

PeerScout

Submission process (simplified)



Also see Labs post: [Peer Review: New initiatives to enhance the value of eLife's process](#)

PeerScout - Overview

BY MANUSCRIPT

Search Type: Reviewer | Manuscript number (last 5 digits): 18449

Legend:

- Main manuscript
- Potential reviewer
- Potential reviewer with review duration
- Corresponding author of selected manuscript
- Early career reviewer
- Related manuscript
- 100 Combined score (keyword & similarity)

Title: "Antiparallel protocadherin homodimers use distinct affinity- and specificity-mediating regions in cadherin repeats 1-4" 10.7554/eLife.18449

Authors: Dr. John Nicoludis, Bennett Vogt, Anna Green, Charlotta Schärfe, Dr. Debora Marks, Prof. Rachelle Gaudet

Senior Editors: [Name]

Subject areas: Computational and Systems Biology, Structural Biology and Molecular Biophysics

Abstract: "Protocadherins (Pcdhs) are cell adhesion and signaling proteins used by neurons to develop and maintain neuronal networks, relying on trans homophilic interactions between their extracellular cadherin (EC) repeat domains. We present the structure of the antiparallel EC1-4 homodimer of human PcdhyB3, a member of the γ subfamily of clustered Pcdhs. Structure and sequence comparisons of α , β , and γ clustered Pcdh isoforms illustrate that subfamilies encode specificity in distinct ways through diversification of loop region structure and composition in EC2 and EC3, which contains isoform-specific conservation of primarily polar residues. In contrast, the EC1/EC4 interface comprises hydrophobic interactions that provide non-selective dimerization affinity. Using sequence coevolution analysis, we found evidence for a similar antiparallel EC1-4 interaction in non-clustered Pcdh families. We thus deduce that the EC1-4 antiparallel homodimer is a general interaction strategy that evolved before the divergence of these distinct protocadherin families."

Review Time: Overall: 5.4 days (avg over 4 reviews), 0 reviews in progress, 0 reviews awaiting response, 0 reviews declined
Last 12 months: n/a

Scores: 100 (max across manuscripts)

Reviewer: [Name] (early career reviewer) ORCID Crossref

Review Time: Overall: 6.6 days (avg over 3 reviews), 0 reviews in progress, 0 reviews awaiting response, 1 review declined
Last 12 months: 7.3 days (avg over 2 reviews), 0 reviews in progress, 0 reviews awaiting response, 1 review declined

Author of: [Name]

Scores: 58 (max across manuscripts)

Suggested reviewers

Dr. [Name] [Q](#)
[ORCID](#) [Crossref](#)

Review Time: Overall: 5.4 days (avg over 4 reviews), 0 reviews in progress, 0 reviews awaiting response, 0 reviews declined
Last 12 months: n/a

Scores: **100** (max across manuscripts)

Dr. [Name] (early career reviewer) [ORCID](#) [Crossref](#)
[ORCID](#) [Crossref](#)

Review Time: Overall: 6.6 days (avg over 3 reviews), 0 reviews in progress, 0 reviews awaiting response, 1 review declined
Last 12 months: 7.3 days (avg over 2 reviews), 0 reviews in progress, 0 reviews awaiting response, 1 review declined

Author of: [Title] [Link]
[Title] [Link]
[Link]

Scores: **58** (max across manuscripts)

Conclusion

- Lessons learned:
 - Work more closely with editors
 - Need for more integrated tools
- Next:
 - Integrate with our new data tools
 - More interpretable results to increase trust

PeerTax

(work by **Alessio Caciagli**)

Reviews can be long.. (e.g. 918 words)

In this manuscript Barry, Behet and colleagues address an important question aiming to possibly appoint an effective role in protection from clinical malaria to liver-stage immunity acquired naturally in malaria endemic areas. I believe the work presented is a significant contribution to our understanding of naturally acquired immunity to pre-erythrocytic stages of *P. falciparum* and I make a few comments and suggestions that may improve the current version of the manuscript.

The data shows that ABs developed during past malaria cases can reduce sporozoite motility and hepatocyte invasion *in vitro* suggesting that ABs acquired during natural infections can reduce new liver-stage infections. However, the contribution of a possible similar effect *in vivo* is, in my view, less clear from the data presented, and more caution may be needed to discuss the results. All the children followed throughout the study became parasite positive by qPCR, indicating that even children with strong inhibiting ABs were unable to block liver-stage infection efficiently. Furthermore, *high response to asexual stage lysate* might be confounding the analysis. I would add to the manuscript (or supplemental) figures the association of CSP ABs with asexual stage lysate ABs, and also the association of whole SPZ ABs with asexual stage lysate ABs, to give the reader an idea of how close these parameters are.

I also suggest including, if that has not been done already, the *Reported bed net use* in the multivariate analyses, as this could also be a factor increasing time to PCR positive and clinical malaria.

The authors cite the study by Tran et al.¹ where it was shown that, in Mali, time to PCR positive was independent of age, while time to clinical malaria increased with age, and where as stated it was concluded that *there no or very limited evidence for an age dependent acquisition of immunity protecting from infection*. Similarly, in the present manuscript, *in vitro* functional data of higher humoral response against pre-erythrocytic stages does not (independently of blood-stage immunity) protect from infection. So, I would rephrase the last sentence in first paragraph of the discussion to add a bit more caution in interpreting what may be causing partial protection.

At the end of section *Evidence of natural risk-modifying pre-erythrocytic immunity* the authors should, in my view, clearly state that *High gliding inhibition activity* does not independently associate in a statistically significant way, with *protection against falciparum infection* on the multivariate analyses where blood stage immunity was included; the P value is above 0.05 (0.055) and the CI includes 1, making the relative risk not statistically significant.

I believe the manuscript could be improved by presenting the quantitative analysis of the 18s qPCR upon first parasite detection and determine if there is a negative association with the inhibitory capacity of the individuals' ABs. It would also be very interesting to question if time from first PCR positive to time of presentation of symptoms is different between poor and strong *in vitro* inhibitors. If the *in vitro* data showing gliding inhibition and reduced hepatocyte invasion are significant *in vivo*, one would expect a lower inoculum in the liver and thus a lower parasitaemia on the first PCR positive time-point. And then potentially a slower progression to clinical malaria. I believe with the data generated in this manuscript these analyses could be done, and would enrich the story.

It is not totally clear to me how individuals were selected for the flow cytometry assays. Survival, gliding inhibition, CSP, LSA1 and asexual lysate ELISAs were performed for the 51 participants, but flow cytometry data presented in fig2 D and E was obtained from 17 Burkinabes only; how were those selected and what is their time to PCR+ in the survival analysis. If they are the 8 poor and 8 strong inhibitors as defined by their gliding inhibition it should be stated in the methods (seems to be so, given supFig3, but there is one extra?). I also suggest to pinpoint these 8 poor and 8 strong inhibitors in fig1 so that the reader would be informed of their time to PCR+ and time to malaria symptoms.

I would be more cautious when citing ref 31, I believe the study by Michael Stewart et al.² shows that non-motile SPZ are unable to invade, but is not clearly showing a direct association between % of human AB affecting motility and those levels correlating directly with invasion either.

Minor points:

In table 1, I would not refer to the 6 children who were PCR positive at the 3 weeks after treatment time-point as *Persisting parasites post-treatment* I do not think that it can be excluded that the children were re-infected after clearance of PQ.

I recommend adding a brief description of the method in ref 45 in the section *In vitro sporozoite infectivity assay of a human hepatoma cell line* in material and methods.

On page 3 below table 1 there is mention to *field PCR* which may be a mistake.

The data from the *in vitro* gliding inhibition by LSA IgG seems to be not shown. I think it should be clarified in the text that that is indeed the case. Likewise, if the LSA-1-specific IgG antibodies correlation with sporozoite invasion inhibition is data not shown I would clearly state it in the text.

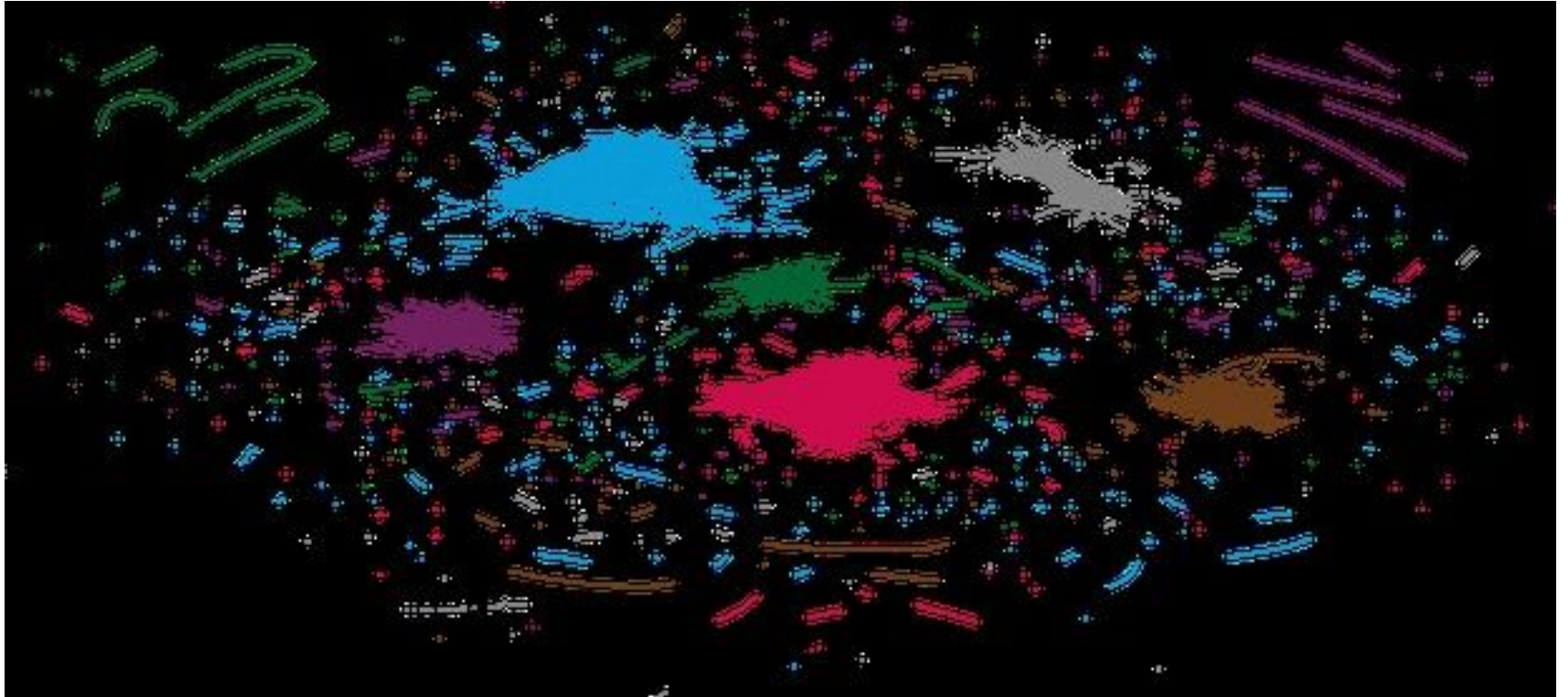
In figS2A I would specify that is IgG in the figure x axis and use the label CSP IgG titer instead of CSP antibody titer.

Figure S3D is called before Figure S3C, I would call figures in ascending and alphabetical order instead.

Main clusters identified

- Figures/Non-textual content
- Stats/Analysis/Techniques
- Impact/Novelty
- Text/Exposition clarity
- Previous literature
- Main discussion

HDBSCAN Output (main topics plus residuals)



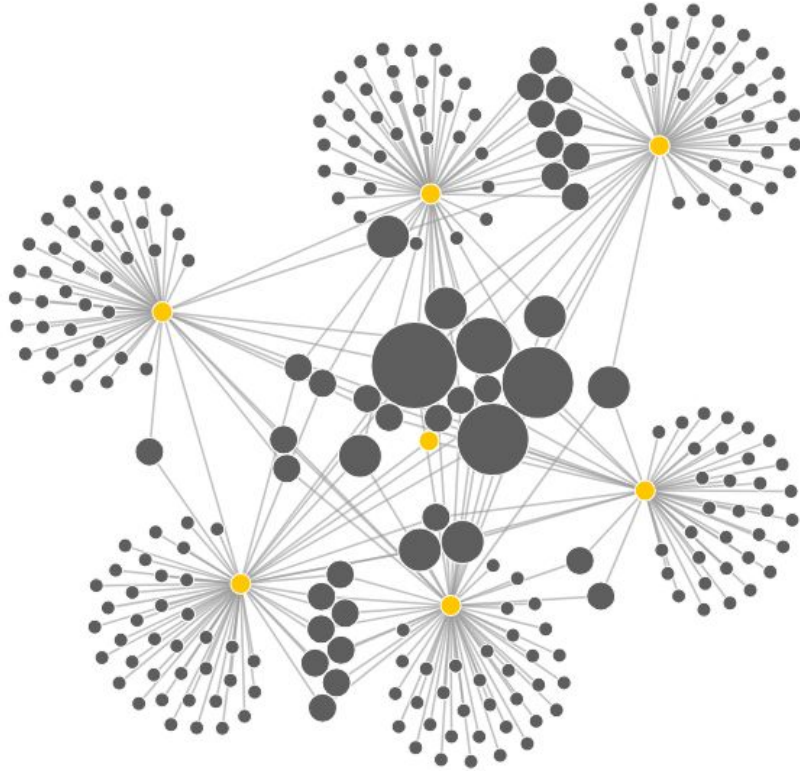
Summary

- Explored automatically structure peer-review content
- Created good first model
- Next:
 - More open training data
 - Better categories
 - More labelled data
 - Supervised model
 - Define use-cases

Citation Sentiment

(work by David Ciudad)

Example citation network visualisation



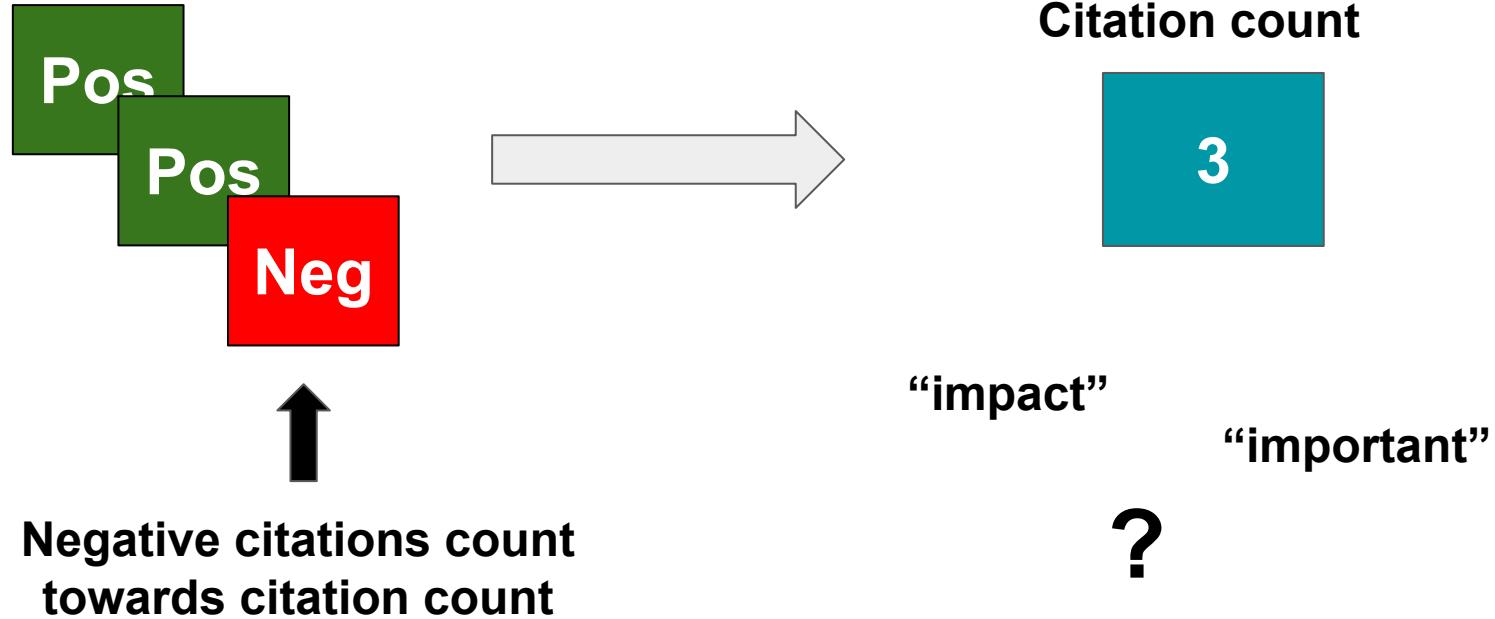
citation gecko 

Example citation

“Our results contrast with the high rate of XMRV detection reported by *Lombard et al.* among both CFS patients and controls, but are in agreement with recent data reported in two large studies in the UK and in the Netherlands...”

([Switzer et al., 2010](#))

Why does the sentiment matter?



Summary

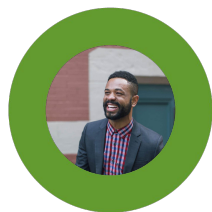
- Explored initial model that is better than off-the-shelf models
- Current datasets not large enough
- Next:
 - Find or create larger dataset
 - Extend to citation function
 - Define good use-cases

Summary

Data Science projects in context



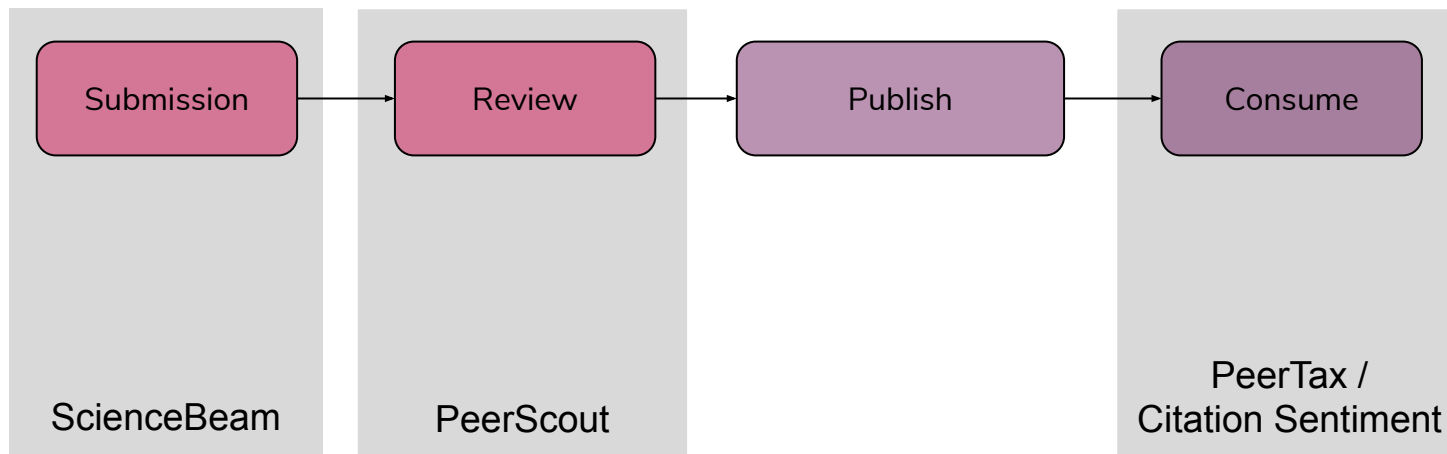
Amy
Author



Hubert
Reviewing Editor / Editorial Staff



Philip
Reader



Resources

- [ScienceBeam Labs blog](#)
- [PeerTax slides](#) (and [labs blog](#))
- [Citation Sentiment slides for longer talk](#) (as part of [Workshop on Open Citations](#))



Thank you

hhmi | Howard Hughes
Medical Institute



*Knut and Alice
Wallenberg
Foundation*

A black and white photograph showing the back of a person's head and shoulders in silhouette. They are looking at a television screen that displays a dense, noisy, grainy pattern. The word "Fin" is overlaid in white text on the person's neck area.

Fin