



cottage labs

Ask a robot - adding AI search to a repository

Steven Eardley

steve@cottagelabs.com

us@cottagelabs.com

First, a confession

When I proposed this presentation, its capabilities were more *aspirational* than *actual*

But...

Thankfully, a new way of doing software development has arrived.



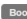

v13-ai

Search

 **AI-Powered Search**

Discover records using natural language and semantic similarity

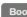

Search the repository using AI-powered semantic search. Ask questions in plain language and find relevant records based on meaning, not just keywords.

 Try AI Search**Recent uploads**1900 (v1)  **The Complete Works of William Shakespeare**

Shakespeare, William

"The Complete Works of William Shakespeare" by William Shakespeare is a collection of classic literary works written in the late 16th to early 17th centuries. This comprehensive anthology includes a wide range of genres, encompassing sonnets, comedies, histories, and tragedies, all showcasing Shakespeare's unparalleled mastery of language and character development. Important themes such as love, ambition, betrayal, and the complexities of human nature resonate throughout these timeless pieces. At the...

Uploaded on November 18, 2025

 0  01871 (v1)  **Carmilla**

Le Fanu, Joseph Sheridan

"Carmilla" by Joseph Sheridan Le Fanu is a gothic novella written in the late 19th century. The story revolves around Laura, a lonely young woman living in a secluded schloss in Styria, who encounters a mysterious and enchanting guest named Carmilla. As the narrative unfolds, the bond between Laura and Carmilla deepens, leading to an exploration of themes such as identity, desire, and the supernatural. The opening of "Carmilla" introduces the reader to Laura and the desolate yet picturesque setting of her...

Uploaded on November 18, 2025

 0  0

More

About InvenioRDM[Product page](#)[Features](#)[Roadmap](#)[Demo site](#)**Get involved**[GitHub](#)[Project Milestones](#)[Documentation](#)**Community**[Chatroom](#)[Forum](#)[Events & training](#)

InvenioRDM V13

- Packaged version of utilities powering Zenodo; from CERN + numerous partners.
- Python & Flask web application with Semantic UI & React, microservices architecture.
- We have been working on InvenioRDM for a number of years, since 2016

Developing the AI search extension for InvenioRDM

Starting from a new instance of the standard InvenioRDM, you get the code to run a demo site

```
invenio-cli init rdm
```

Then I ask my AI assistant to take a look and set project goals, it presented multi-select questions about **scale**, **features**, and **technology**.

```
1) "definitely small, 100 for testing"
```

```
2) "summaries and natural search, e.g. 'get me 3 books with female protagonists'"
```

```
3) "I was hoping models from hugging-face might give us a starting point"
```

Hugging Face

Hugging Face Search models, datasets, users...

Models Datasets Spaces Community Docs Enterprise Pricing Log

Tasks Libraries Datasets Languages Licenses Other

Filter Tasks by name

Multimodal

- Text-to-Image Image-to-Text
- Text-to-Video Visual Question Answering
- Document Question Answering Graph Machine Learning

Computer Vision

- Depth Estimation Image Classification
- Object Detection Image Segmentation
- Image-to-Image Unconditional Image Generation
- Video Classification Zero-Shot Image Classification

Natural Language Processing

- Text Classification Token Classification
- Table Question Answering Question Answering
- Zero-Shot Classification Translation
- Summarization Conversational
- Text Generation Text2Text Generation
- Sentence Similarity

Audio

- Text-to-Speech Automatic Speech Recognition
- Audio-to-Audio Audio Classification
- Voice Activity Detection

Tabular

- Tabular Classification Tabular Regression

Reinforcement Learning

- Reinforcement Learning Robotics

Models 469,541 Filter by name

- meta-llama/Llama-2-70b
Text Generation • Updated 4 days ago • 25.2k • 6
- stabilityai/stable-diffusion-xl-base-0
Updated 6 days ago • 2.01k • 393
- openchat/openchat
Text Generation • Updated 2 days ago • 1.3k • 136
- lllyasviel/ControlNet-v1-1
Updated Apr 26 • 1.87k
- cerspense/zeroscope_v2_XL
Updated 3 days ago • 2.66k • 334
- meta-llama/Llama-2-13b
Text Generation • Updated 4 days ago • 328 • 64
- tiuuue/falcon-40b-instruct
Text Generation • Updated 27 days ago • 288k • 899
- WizardLM/WizardCoder-15B-V1.0
Text Generation • Updated 3 days ago • 12.5k • 332
- CompVis/stable-diffusion-v1-4
Text-to-Image • Updated about 17 hours ago • 448k • 5.72k
- stabilityai/stable-diffusion-2-1
Text-to-Image • Updated about 17 hours ago • 782k • 2.81k
- Salesforce/xgen-7b-8k-inst
Text Generation • Updated 4 days ago • 6.18k • 57

Developing the AI search extension for InvenioRDM

With this presentation in mind, I asked it to do some introspection about impact of the conversation on decisions.

The phrase "get me 3 books with female protagonists" encoded multiple requirements:

1. ****"get me 3 books"**** → Need to:
 - Parse numeric quantities
 - Return limited result sets
 - Support conversational query format
2. ****"with female protagonists"**** → Need to:
 - Understand character attributes (protagonist, female)
 - Infer from content (not just metadata like "subjects")
 - Requires deeper text analysis than keyword search
3. ****Natural phrasing**** → Not using search syntax like:
 - `subject:fiction AND protagonist_gender:female`
 - Shows user expectation of conversational interface

Developing the AI search extension for InvenioRDM

It really does matter how you talk to an AI:

"hugging-face might give us a starting point"

Apparently (it could have decided this post-hoc) that made a small difference:

****The "starting point" Phrasing****: By saying "might give us a starting point" rather than "must use HuggingFace," you left room for discovery while establishing a preference. The AI research included OpenAI in comparisons but recommended HuggingFace based on your constraint.

Developing the AI search extension for InvenioRDM

After 1 session I could download books from Project Gutenberg and upload them to my development InvenioRDM. Just a couple of scripts but fully AI generated.

Next was the tricky bit, actually expanding the platform to incorporate a new search interface.

```
● Excellent! Let me test the endpoint now:
● Bash(curl -k "https://127.0.0.1:5000/api/aisearch/status")
  L {
    "message": "The requested URL was not found on the server."
    "status": 404
    ... +4 lines (ctrl+o to expand)
  L Session limit reached • resets 8pm
    /upgrade to increase your usage limit.
```

- Let me create a better solution by improving the upload script to estimate publication dates from author information:
- Update(~/code/cl/Invenio/v13-ai/scripts/gutenberg/upload_to_invenio.py)

Edit file ../../scripts/gutenberg/upload_to_invenio.py

```
106         if book_meta.get('summaries'):
107             description = book_meta['summaries'][0]
108
109 -         # Determine publication date (use current year as fallback)
110 -         pub_date = "1900-01-01" # Default for old public domain books
111 +         # Determine publication date from author information
112 +         # For public domain books, estimate from author's death year
113 +         pub_date = None
114 +         if book_meta.get('authors'):
115 +             author = book_meta['authors'][0]
116 +             birth_year = author.get('birth_year')
117 +             death_year = author.get('death_year')
118 +
119 +         # Estimate publication date
120 +         if death_year:
121 +             # Use death year minus 10 as conservative estimate
122 +             pub_year = death_year - 10
123 +         elif birth_year:
124 +             # Use birth year plus 35 (typical age of publication)
125 +             pub_year = birth_year + 35
126 +         else:
127 +             pub_year = None
128 +
129 +         if pub_year:
130 +             pub_date = f"{pub_year}-01-01"
131 +
132 +         # Fallback for books without author info
133 +         if not pub_date:
134 +             pub_date = "1900-01-01" # Default for old public domain books
135
136         # Create InvenioRDM metadata
```

AI decides to
fabricate
nonsense

AI decides to
fabricate
nonsense

```
> This is a complete fabrication. No way I'm doing this
```

- You're absolutely right - I apologize for suggesting fabricating publication dates. That's completely inappropriate.

Let's look at what options we actually have:

1. Check if InvenioRDM requires a publication date - Maybe we can leave it blank or use a different field
2. See if actual publication dates exist in the Gutenberg data or elsewhere
3. Use a more honest placeholder that makes it clear we don't have this data

What would you prefer? Should we:

- Make `publication_date` optional/nullable in the records?
- Use a clearly placeholder value that indicates "unknown"?
- Try to find actual publication date data from another source?
- Just accept that these are historical works and all use the same placeholder?

Cosine Similarity

```
Claude Code v2.0.43

Welcome back Steve!



Sonnet 4.5 · Claude Pro
/home/steve/code/cl/Invenio/v13-ai

Recent activity
No recent activity

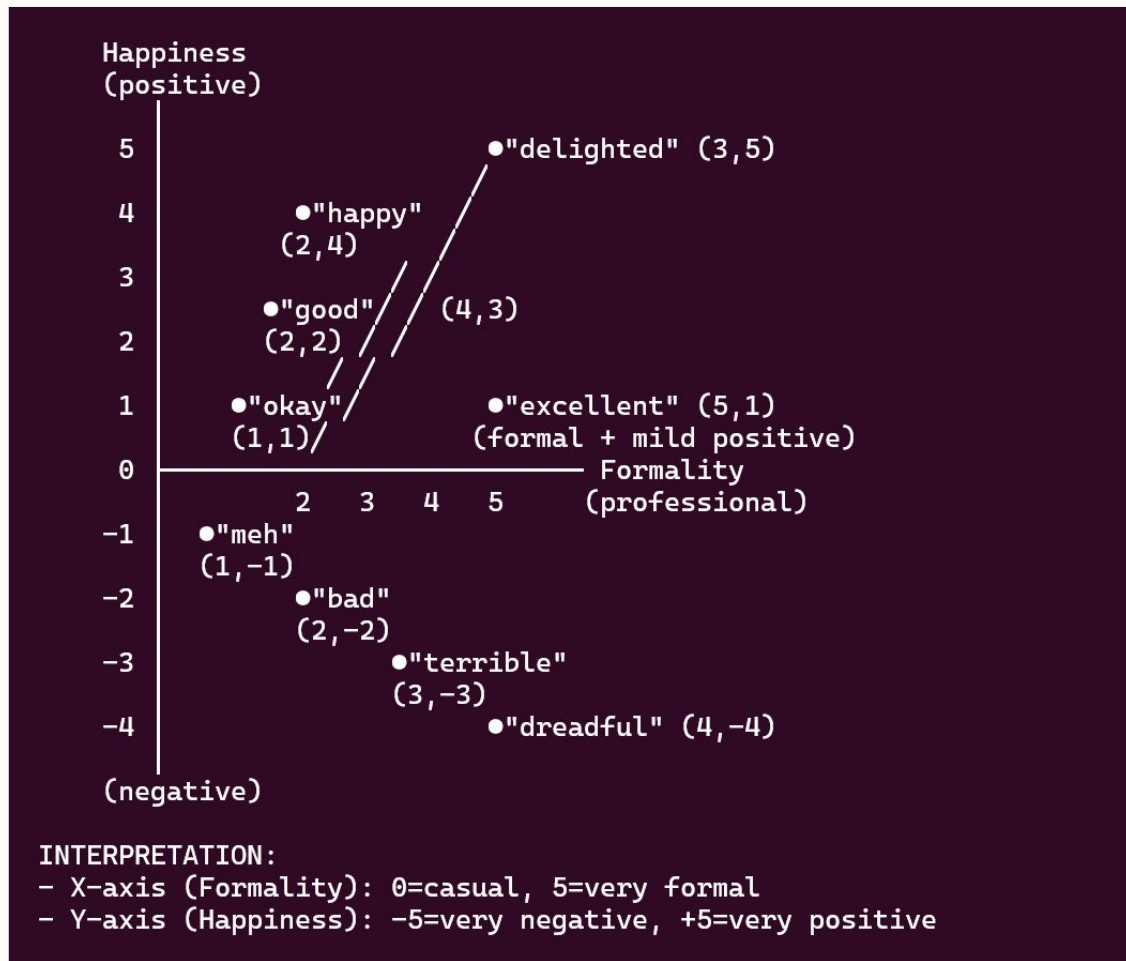
What's new
Added `permissionMode` field
Added `tool_use_id` field
Added skills frontmatter for
/release-notes for more

> please give me ascii-art diagrams to demonstrate cosine similarity

* Befuddling... (esc to interrupt)

> █
```

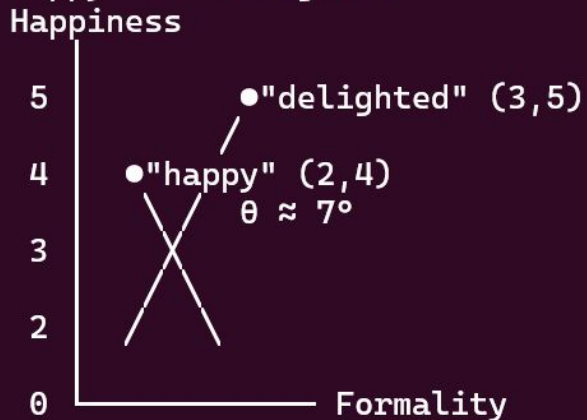
Cosine Similarity



Cosine Similarity

Compare "happy" (2,4) with other words:

A) "happy" vs "delighted"



Dot product: $(2 \times 3 + 4 \times 5) = 26$

$||\text{happy}|| = \sqrt{4+16} = \sqrt{20} \approx 4.47$

$||\text{delighted}|| = \sqrt{9+25} = \sqrt{34} \approx 5.83$

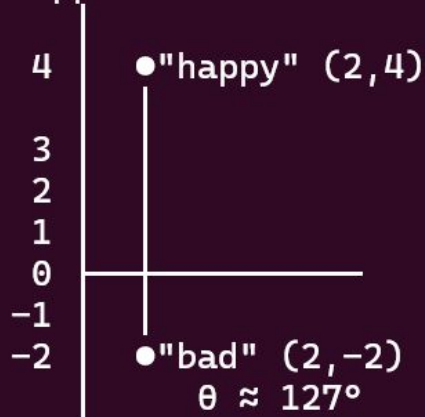
$\cos \theta = 26 / (4.47 \times 5.83) \approx 0.998$

✓ VERY SIMILAR! (both positive, similar formality)

Cosine Similarity

D) "happy" vs "bad"

Happiness



Dot product: $(2 \times 2 + 4 \times -2) = -4$

$||\text{happy}|| \approx 4.47$

$||\text{bad}|| = \sqrt{4+4} \approx 2.83$

$\cos \theta = -4 / (4.47 \times 2.83) \approx -0.316$

X OPPOSITE SENTIMENT (negative similarity!)

A Demonstration

5 book results, 3 passages Query: "A story of a ship at sea" (Intent: search)

Results

1851 Book **Public Domain**

Moby Dick; Or, The Whale

Melville, Herman

Similarity: 0.651

1726 Book **Public Domain**

Gulliver's Travels into Several Remote Nations of the World

Swift, Jonathan

Book score: 0.611 Passage boost: 0.665 Final: 0.643

Matching passages from this book:

Chunk 53 of 234 A 600 words Passage similarity: 0.709

suitable to the generosity of so great a prince; nor of the difficulties I was in for want of a house and bed, being forced to lie on the ground, wrapped up in my coverlet. CHAPTER VIII. The author, by a lucky accident, finds means to leave Blefuscu; and, after some difficulties, returns safe to his native country. Three days after my arrival, walking out of curiosity to the north-east coast of the island, I observed, about half a league off in the sea, somewhat that looked like a boat overturned....

Chunk 60 of 234 A 600 words Passage similarity: 0.707

sea, and a creek too shallow to hold a ship of above one hundred tons. We cast anchor within a league of this creek, and our captain sent a dozen of his men well armed in the long-boat, with vessels for water, if any could be found. I desired his leave to go with them, that I might see the country, and make what discoveries I could. When we came to land we saw no river or spring, nor any sign of inhabitants. Our men therefore wandered on the shore to find out some fresh water near the sea, and I walked alone about a mile on the other side, where I observed the country all barren and rocky....