

PathoGeneMap

*Empower biomedical researchers with
AI-enhanced exploration of disease causing
genes*

Wesley Chang, Almicia Dunson,
Max Kaufmann, Christian Lee, Andrew McCall

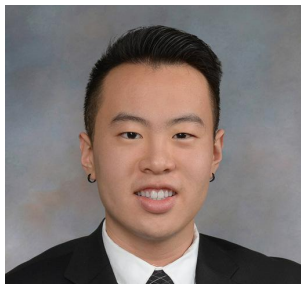


OUR TEAM



Almicia Dunson

PM & Web Dev



Wesley Chang

Infrastructure & Web
Dev



Max Kaufmann

SME & Data
Engineering



Christian Lee

ML Engineering



Andrew McCall

Infrastructure & Data
Engineering

THE PROBLEM



- Finding new treatment targets for human diseases
- More than 20,000 genes to choose from
- Which ones should I target?

THE PROBLEM



- Standing on the shoulders of giants

PubMed[®]

- Provides > 23 million biomedical literature abstracts

THE PROBLEM

Experiments in PubMed are described in **natural language only!**

Intracerebral infection of mice with mouse hepatitis virus, a member of the Coronaviridae family, reproducibly results in an acute encephalomyelitis that progresses to a chronic demyelinating disease. The ensuing neuropathology during the chronic stage of disease is primarily immune mediated and similar to that of the human demyelinating disease multiple sclerosis. Secretion of chemokines within the CNS signals the infiltration of leukocytes, which results in destruction of white matter and neurological impairment. The CC chemokine ligand (CCL5) is localized in white matter tracts undergoing demyelination, suggesting that this chemokine participates in the pathogenesis of disease by attracting inflammatory cells into the CNS. In this study, we administer a mAb directed against CCL5 to mice with established mouse hepatitis virus-induced demyelination and impaired motor skills. Anti-CCL5 treatment decreased T cell accumulation within the CNS based, in part, on viral Ag specificity, indicating the ability to differentially target select populations of T cells. In addition, administration of anti-CCL5 improved neurological function and significantly ($p < \text{or} = 0.005$) reduced the severity of demyelination and macrophage accumulation within the CNS. These results demonstrate that the severity of CNS disease can be reduced through the use of a neutralizing mAb directed against CCL5 in a viral model of demyelination.

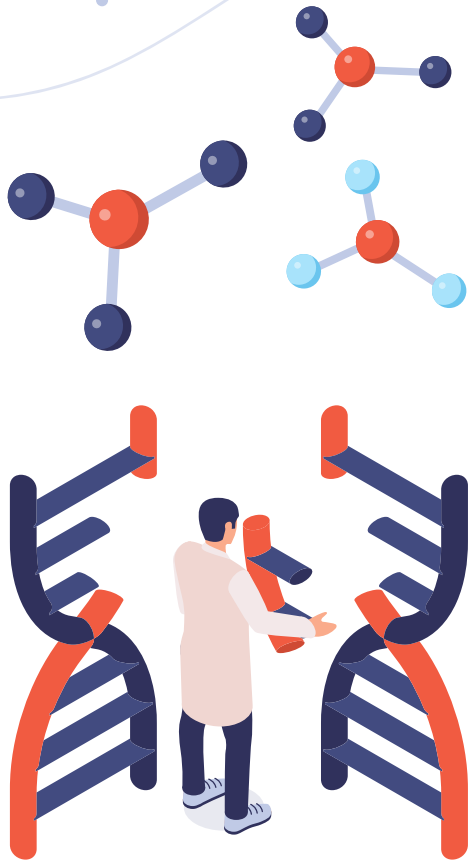


Inhibition of CCL5 ameliorated a multiple sclerosis disease model



**THE STATE OF
THE ART IS ...**

... READING



OUR MISSION

Empower biomedical researchers with AI-enhanced exploration of disease causing genes.

THE IMPACT

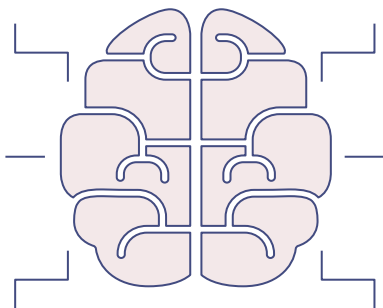
Number of genes interested in:	100
Amount of articles per gene:	900
Avg. minutes to determine article relevance:	2
Total days to find articles(8hr day):	375

A **couple of minutes** for a researcher to find relevant articles

375 Days

- Time for a researcher without our product
- Time for a researcher *with* our product

WHY HAS IT **NOT** BEEN DONE?



There is **no** model.



There is **no** labeled training data.



There is **no** funding.

THE TASK

Positive example

Intracerebral infection of mice with mouse hepatitis virus, a member of the Coronaviridae family, reproducibly results in an acute encephalomyelitis that progresses to a chronic demyelinating disease. The ensuing neuropathology during the chronic stage of disease is primarily immune mediated and similar to that of the human demyelinating disease multiple sclerosis. Secretion of chemokines within the CNS signals the infiltration of leukocytes, which results in destruction of white matter and neurological impairment. The CC chemokine ligand (CCL5) is localized in white matter tracts undergoing demyelination, suggesting that this chemokine participates in the pathogenesis of disease by attracting inflammatory cells into the CNS. In this study, we administer a mAb directed against CCL5 to mice with established mouse hepatitis virus-induced demyelination and impaired motor skills. Anti-CCL5 treatment decreased T cell accumulation within the CNS based, in part, on viral Ag specificity, indicating the ability to differentially target select populations of T cells. In addition, administration of anti-CCL5 improved neurological function and significantly ($p < \text{or} = 0.005$) reduced the severity of demyelination and macrophage accumulation within the CNS. These results demonstrate that the severity of CNS disease can be reduced through the use of a neutralizing mAb directed against CCL5 in a viral model of demyelination.

Model
input

Inhibition of CCL5 ameliorated a
multiple sclerosis disease model

Perturbation type	Inhibition
Gene	CCL5
Effect	Amelioration
Disease model	Multiple sclerosis

Model
output

Database

Negative example

Experimental autoimmune encephalomyelitis (EAE) is a classical experimental model of multiple sclerosis (MS), an autoimmune disease of the central nervous system (CNS). Previous reports have suggested that matrine (MAT), a quinolizidine alkaloid derived from the herb Radix Sophorae Flave, could inhibit clinical EAE, but its mechanism of action is not clear. Our present study showed that MAT treatment resulted in dose-dependent reduction in neurological scores. Consistent with this observation, infiltration of inflammatory cells and demyelination in the CNS were also significantly suppressed. We further studied the mechanism underlying these effects of MAT by determining whether this treatment influences expression of molecules that are involved in the activation and migration of inflammatory cells. Our results showed that MAT significantly inhibited expression and production in the CNS of ICAM-1 and VCAM-1, key adhesive molecules, and CCL3 and CCL5, key chemokines, that attract inflammatory cells into the CNS. Furthermore, the TLR4/MD2 pathway, which plays an important role in the induction of Th1/Th17 cells in EAE, was also significantly inhibited. Together, our study not only demonstrates that MAT may be a novel therapeutic option for the treatment for MS, but also provides further information on mechanisms underlying the effect of MAT treatment.

No specific gene was perturbed in vivo.



THE TASK

Positive example

Intracerebral infection of mice with mouse hepatitis virus, a member of the Coronaviridae family, reproducibly results in an acute encephalomyelitis that progresses to a chronic demyelinating disease. The ensuing neuropathology during the chronic stage of disease is primarily immune mediated and similar to that of the human demyelinating disease multiple sclerosis. Secretion of chemokines within the CNS signals the infiltration of leukocytes, which results in destruction of white matter and neurological impairment. The CC chemokine ligand (CCL5) is localized in white matter tracts undergoing demyelination, suggesting that this chemokine participates in the pathogenesis of disease by attracting inflammatory cells into the CNS. In this study, we administer a mAb directed against CCL5 to mice with established mouse hepatitis virus-induced demyelination and impaired motor skills. Anti-CCL5 treatment decreased T cell accumulation within the CNS based, in part, on viral titer, suggesting the ability to differentially target select populations of T cells. In addition, administration of anti-CCL5 improved neurologic function and significantly ($p < 0.05$) reduced the severity of demyelination and axonal loss. These results demonstrate that the severity of disease can be reduced through inhibition of a neutralizing antigen. Administration of anti-CCL5 in a viral model of demyelination.

Inhibition of CCL5 ameliorated a multiple sclerosis disease model

Perturbation type	Inhibition
Gene	CCL5
Effect	Amelioration
Disease model	Multiple sclerosis

Negative example

Experimental autoimmune encephalomyelitis (EAE) is a classical experimental model of multiple sclerosis (MS), an autoimmune disease of the central nervous system (CNS). Previous reports have suggested that matrine (MAT), a quinolizidine alkaloid derived from the herb Radix Sophorae Flavae, could inhibit clinical EAE, but its mechanism of action is not clear. Our present study showed that MAT treatment resulted in dose-dependent reduction in neurological scores. Consistent with this observation, infiltration of inflammatory cells and demyelination in the CNS were also significantly suppressed. We further studied the mechanism underlying these effects of MAT by determining whether this treatment influences expression of molecules that are involved in the activation and migration of inflammatory cells. Our results showed that MAT significantly inhibited expression and production in the CNS of ICAM-1 and VCAM-1, key adhesive molecules, and CCL3 and CCL5, key chemokines, that attract inflammatory cells into the CNS. Furthermore, TLR4/MD2 signaling, which plays a major role in the induction of Th1/Th17 cells in EAE, was also significantly inhibited. Together, our study not only demonstrates that MAT may be a novel therapeutic option for the treatment of MS, but also provides a potential mechanism underlying the effects of MAT treatment.

No specific gene was perturbed in vivo.

X

Model input

Model output

Database

HOW DOES GPT-4 PERFORM?

Test Data F1 Score

0.84

Cost to run on 23 Million Abstracts?



OUR MODELING APPROACH

PRE-TRAINED MODEL

Sci-Five

- T5-Large
- BioMedical Text
- Encoder-Decoder
- 770M parameters



FINE TUNING

QLoRA

- Efficiency
- Faster Training Time
- Less GPUs needed

FINE TUNING TRAINING DATASET

- Own data set of **892** manually labelled abstracts for training
 - 50% positive examples + 50% matched negative examples
- Public release for the benefit of the community

The screenshot shows a GitHub repository page for 'wesleywchang / 2024-capstone-gene-nlp'. The repository path is 'main / 2024-capstone-gene-nlp / fine_tuning_dataset / fine_tuning_dataset_v1.2b'. A commit by 'MaxKman' is highlighted, with a commit message 'same train / val / test split as for v 1.1' and a commit date of 'yesterday'. Below the commit history is the README content for 'PathoGeneMap finetuning training data v 1.2b', which includes the author's name 'Max Kaufmann' and the date '2024-03-11'. The README also contains a 'Version history' section with two entries: 'v 1.0 2024-02-11: basic version of the dataset' and 'v 1.1 2024-02-21: Added labelled real-life test set of 1000 randomly drawn abstracts from PubMed'.

Name	Last commit message	Last commit date
..		
finetuning_training_data	same train / val / test split as for v 1.1	yesterday
README.md	same train / val / test split as for v 1.1	yesterday
build_training_dataset.Rmd	same train / val / test split as for v 1.1	yesterday
fine_tuning_dataset_v1.2b.Rproj	same train / val / test split as for v 1.1	yesterday
retrieval_of_additional_training_data...	same train / val / test split as for v 1.1	yesterday

PathoGeneMap finetuning training data v 1.2b

Author: Max Kaufmann
Date: 2024-03-11

Version history

- v 1.0 2024-02-11: basic version of the dataset
- v 1.1 2024-02-21: Added labelled real-life test set of 1000 randomly drawn abstracts from PubMed

GPT4

VS

OUR MODEL

Test Data F1 Score

0.84

Cost to run on 23 Million Abstracts



Test Data F1 Score

0.94

SCALING TO ALL OF PUBMED

Hardware

g5.xlarge
g4dn.xlarge
p3.2xlarge
⋮
⋮
⋮

Model

Quantization
LoRA
Merged Weights
⋮
⋮
⋮

Parameters

Input Token Length
Temperature
Batch Size
⋮
⋮
⋮

**Minimize while
maintaining
performance**

**Estimated
Cost**

SCALING TO **ALL OF PUBMED**

Hardware

Model

Parameters

Minimize while

Completed inferences in
parallel, reducing time from
41 days to 5 days

g5
g4c
p3

GPT4

VS

OUR MODEL

Test Data F1 Score

0.84

Cost to run on 23 Million Abstracts?



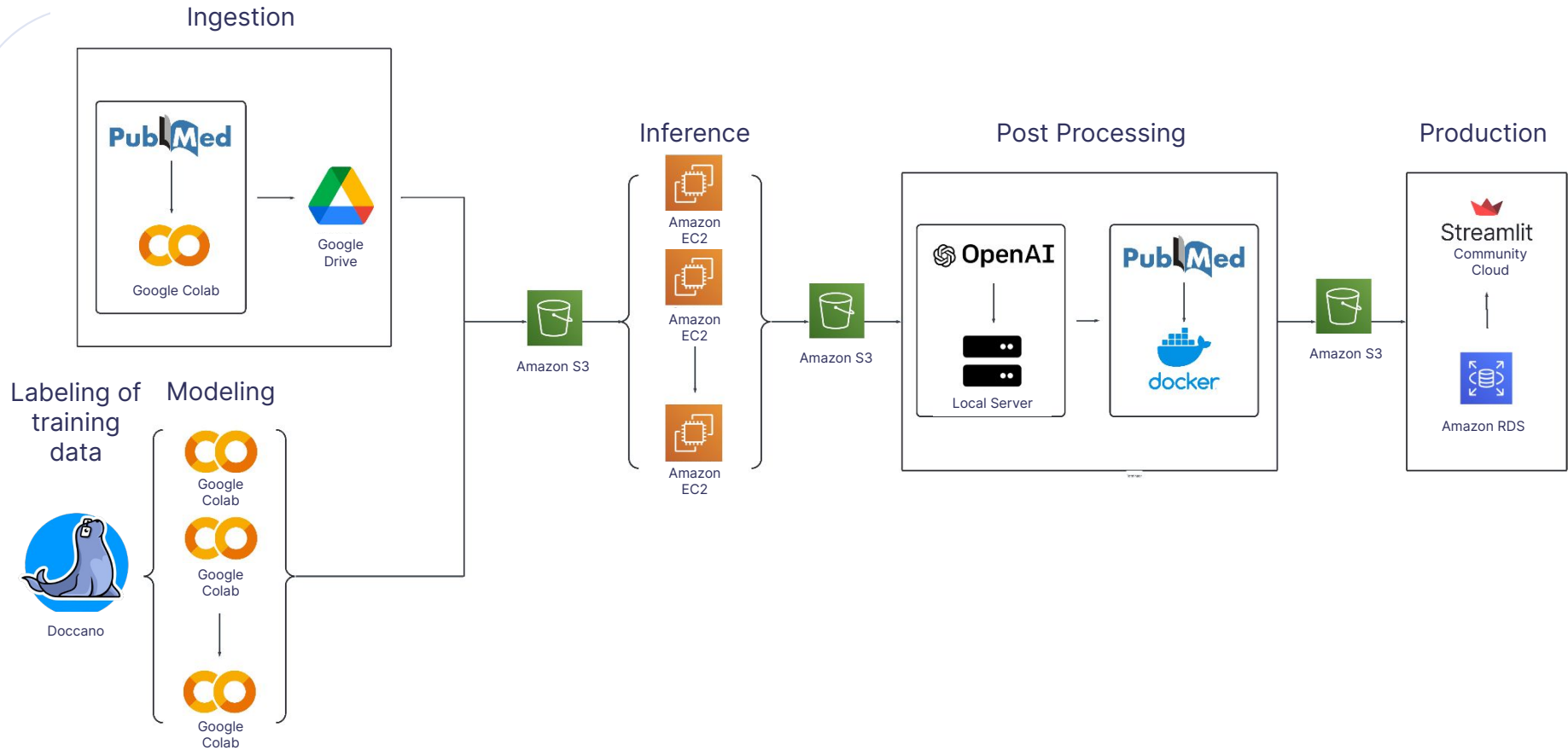
Test Data F1 Score

0.94

Cost to run on 23 Million Abstracts?



HOW DOES IT ALL FIT TOGETHER?

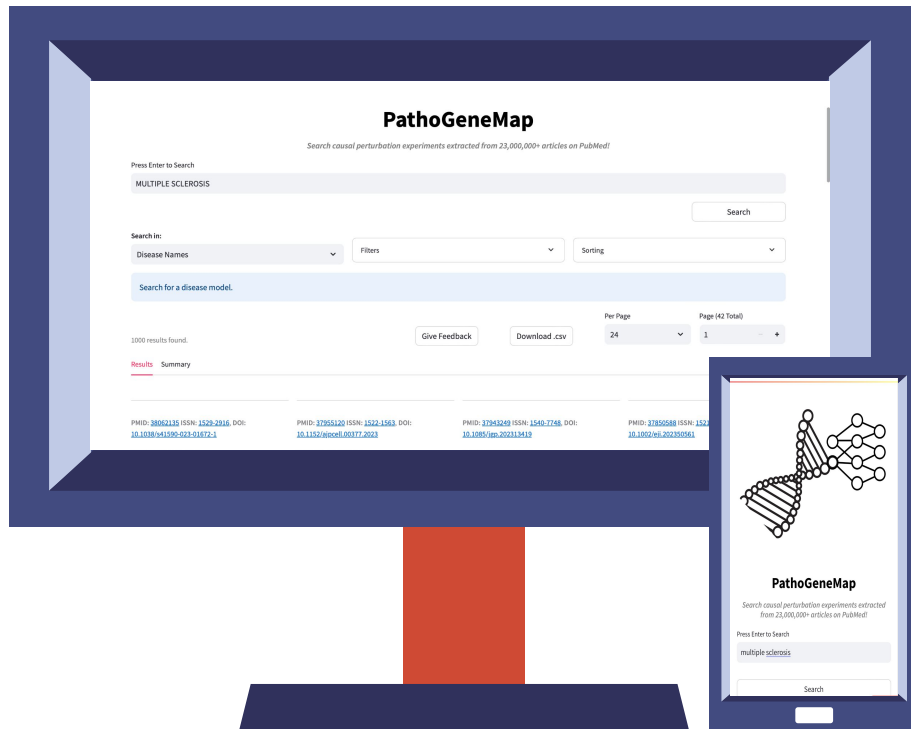


The background features a white space with scattered light blue dots of varying sizes. Three molecular models are positioned around the central text: one in the top right with a red central atom and two light blue peripheral atoms; one in the middle left with a red central atom and two dark blue peripheral atoms; and one in the bottom right with a red central atom and two dark blue peripheral atoms. Faint, light blue curved lines are also visible on the left side of the image.

DEMO

WE ARE LIVE!

VISIT
PATHOGENEMAP.COM
AND TRY OUR APP



REACTIONS FROM OUR USERS



JAN-BRODER ENGLER

MD-PhD
Principal Investigator
University of Hamburg
Germany

“ PathoGeneMap is an excellent starting point for conducting literature research on your gene of interest.

”



KATE ATTFIELD

PhD
Principal Investigator
University of Oxford
United Kingdom

“ A very easy to use tool, which is simple to navigate and gives immediate insight into potential gene functions that are placed directly in the relevant context ”

PRODUCT ROADMAP

WE ARE HERE



Web Portal
Publicly
Available

Get user
feedback and
make updates

Create more
data

Re-train with
new data

Train with
more capable
model

Continuous Tasks - projected running cost of \$25 per month

Update
database

Maintain
website

Retrain
annually with
latest articles

-  In Development
-  Planned for future

CONCLUSION



MISSION

Empower biomedical researchers with AI tools that help determine which genes cause diseases.



APPROACH

We fine tuned LLM's to read abstracts and extract gene information from them, then stored that information in a database that is accessible through a web portal.