



# Master's Research and Projects: FNLCR and Columbia University Partnership

# Master's Research and Projects: FNLCR and Columbia University Partnership

**2:00 PM: Introduction Prof. Michael Robbins**

**2:10 PM: Opening remarks Dr. Ethan Dmitrovsky**

**2:20 PM: Student Presentations**

- Cloud Deployment, Optimization Strategies for Teaching, Training and Collaborative Reproducible Research
- Survey to Identify Emerging Infectious Disease Datasets for Machine Learning
- Survey to Identify Cancer Datasets for Machine Learning
- Q & A

**3:20 PM: Closing remarks Dr. Eric Stahlberg**



# Project Team



Mahitha Kotipalli



Jim Hu



Niranjana Moleyar



Malin Ortenblad



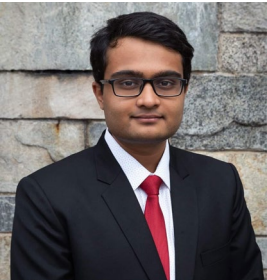
Kerry Hu



Jie Chen



Mengyao He



Om Vaghasia



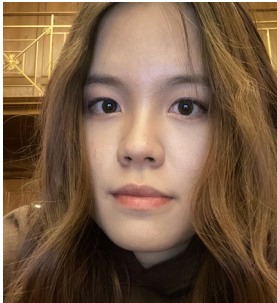
Panagiotis Misirlis



Jiaxi Zhou



Xinyao Wang



Qinwei Zhang



Yue Hu



Zihui Zhou  
Frederick National Laboratory for Cancer Research



Naomi Ohashi, MPA, PMP, ITIL  
Biomedical Informatics and Data  
Science (BIDS)  
Frederick National Laboratory for  
Cancer Research



Ravichandran Sarangan, PhD, PMP  
Biomedical Informatics and Data  
Science (BIDS)  
Frederick National Laboratory for  
Cancer Research



Michael Robbins  
Professor  
Columbia University



Nicole Soder  
TA, Project Manager  
Columbia University





# Survey to Identify Emerging Infectious Disease Datasets for ML

Team Members: Mengyao He, Om Vaghasia, Panagiotis Misirlis

Sep. 8, 2020

# Project Team



Mengyao He  
Columbia University  
MSOR student  
Project 2 - COVID-19  
Prior experience in data analytics and machine learning. Worked with Q-squared to predict futures contracts price movement by Machine Learning models.



Om Vaghasia  
Columbia University  
MSMSE Student  
Project 2 - COVID-19  
Prior experience in data analysis and Natural Language Processing. Worked with International American Supermarkets Corps to optimize their operations.



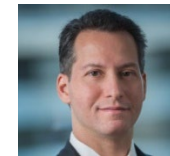
Panagiotis Misirlis  
Columbia University  
MSMSE Student  
Project 2 - SARS  
Prior experience in data analysis, machine learning and deep learning. Worked with a finance firm to create fraud detection tool using ML and with a Non-Profit Organization to create a digital content classifier that is able to detect and count number of Humans in a given picture.



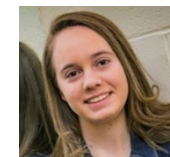
Naomi Ohashi, MPA, PMP, ITIL  
Biomedical Informatics and Data Science (BIDS)  
Frederick National Laboratory for Cancer Research



Ravichandran Sarangan, PhD, PMP  
Biomedical Informatics and Data Science (BIDS)  
Frederick National Laboratory for Cancer Research



Michael Robbins  
Professor  
Columbia University



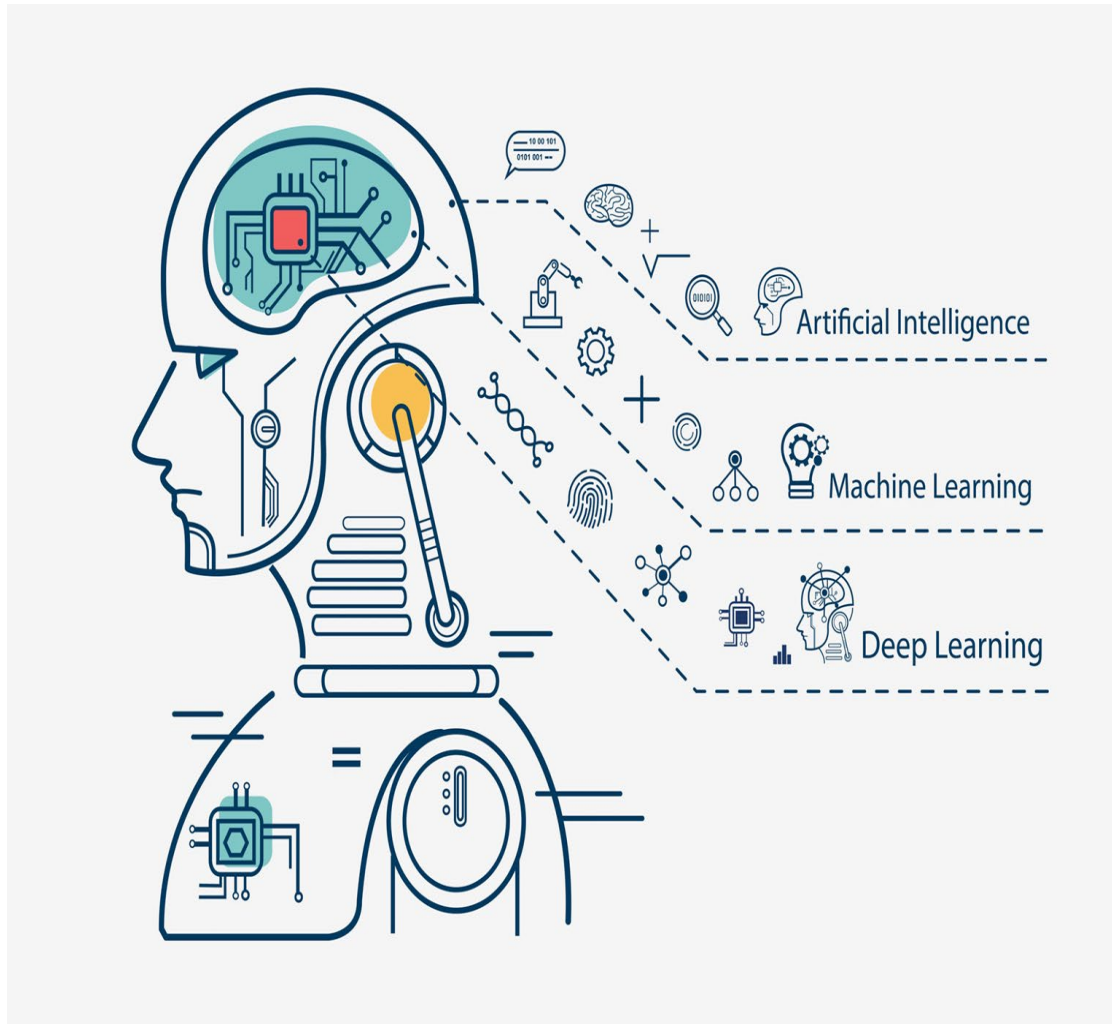
Nicole Lynn Soder  
TA, Project Manager  
Columbia University

# Table of Content

- Project Overview
  - Background
  - Description
  - Goals
- SARS
  - Introduction
  - Methods to Conduct Literature Reviews
  - Results
- COVID-19
  - Introduction
  - Methods to Conduct Literature Reviews
  - Results
- Conclusions
- Future Work



# Project Overview: Background



**Machine-learning (ML)** has become a crucial tool in cancer/biomedical research. The important ingredient for a data science project is the data and despite tremendous progress in ML algorithms, data availability remains a major obstacle.

# Project Overview: Description

- Used PubMed Advanced Search Builder to **collect** and **report** published Machine Learning papers from the area of emerging infectious diseases such as **SARS** (2003) and **COVID-19** (2019)
- For this project, we would prefer the students to focus on ML/AI/Deep-learning modelling efforts that used either **structured** (*Ex. drug SMILES data, DNA/protein, protein-drug binding data*) or **unstructured** (*image data such as X-ray or histopathology etc.*).
- The student(s) will **document the URL location** of the data that accompanies manuscript(s) and **download the raw (uncleaned) data for preparation and analysis.**
- We **also document datasets that are available only on the web** and not associated with any publication



# Project Overview: Goals

- This database will serve as the literature-search starting point for data science projects in NCI/NIH/FNLCR
- The datasets that accompany manuscripts might be ready for modeling and can serve as inputs for proof-of-concept NIH/FNLCR/NCI projects
- Identify any potential trends in the literature in terms of:
  - a) programming languages
  - b) libraries and packages
  - c) specific Machine Learning methodologies (i.e. SVM, Random Forest)

# SARS (2003)

---

Panagiotis Misirlis

# Introduction

## Literature Review:

- The research for SARS-2003 epidemic was mainly focused on the ***spread*** and ***infection rate*** of the virus in different scenarios/cities.
- Another part of the research papers is focused on ***bio-mapping*** the SARS-COV virus.
- The techniques that have been mostly used are :
  - Neural Networks
  - Decision Trees
  - Stochastic Dynamic Models
- They use ***time-series data*** of ***infections/deaths/recoveries*** that were collected in different cities around the world.
- There are certain limitations with SARS literature review and the most significant one is ***time elapsed*** since that epidemic happened. As we are going to see also later on this has a significant impact on the type of techniques that were used.



# Methods to Conduct Literature Reviews

## PubMed Keywords

- Pubmed search engine:
  - <https://pubmed.ncbi.nlm.nih.gov/>
- SARS:
  - SARS 2003 Machine Learning
  - SARS 2003 Forecasting
  - SARS 2003 Prediction
  - SARS 2003 Classification

## Other Methods

- Identified other useful publications from footnotes and references.
- Kaggle
- Github Search

# Results: Major Groups in SARS Studies

## Forecasting & Prediction

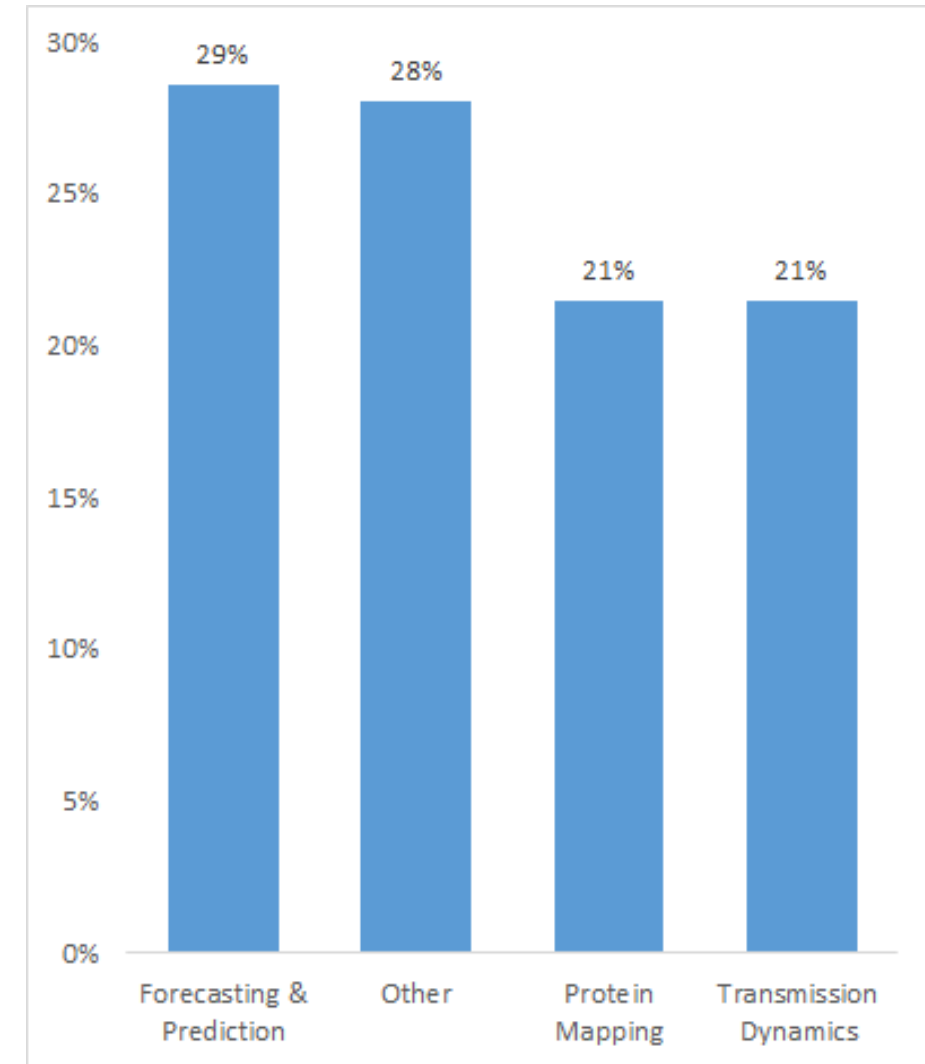
- Forecast the trajectory of the outbreak and quantify the risk of death.
- Predict the transmission rates of the virus under different circumstances

## Protein Mapping

- Prediction Rule generation for SARS-CoV protease cleavage sites.
- Mining SARS-CoV protease cleavage data

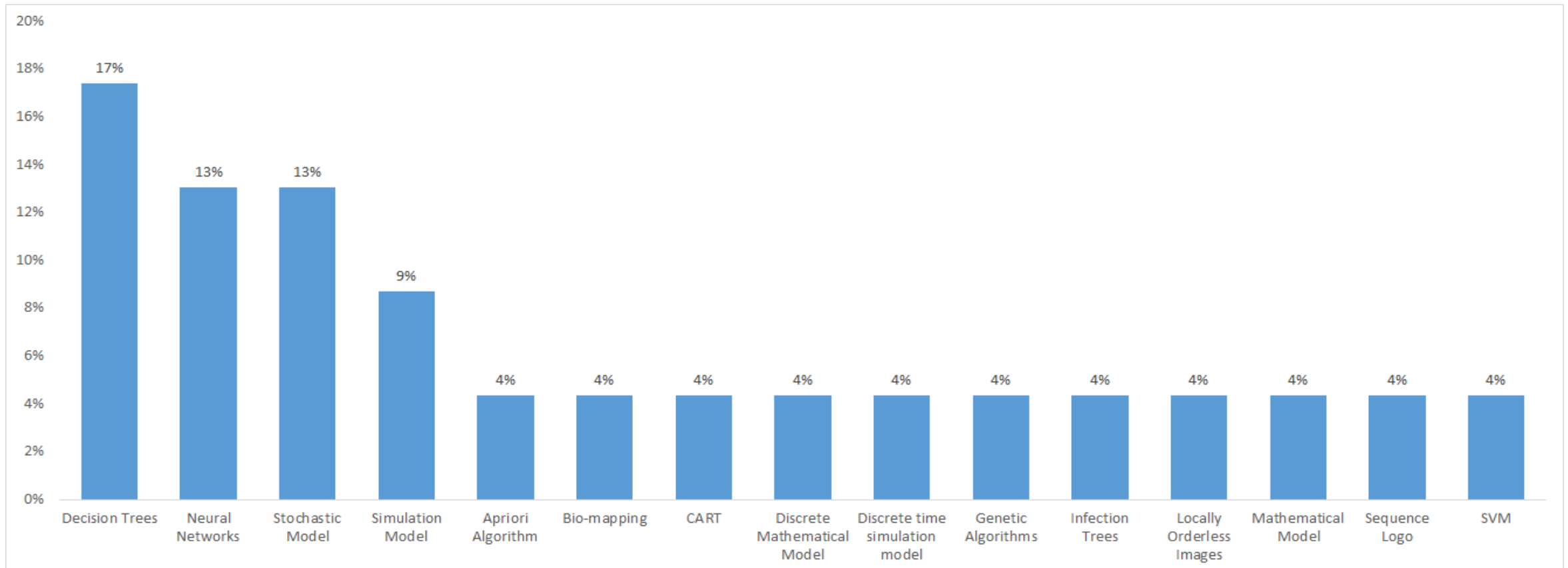
## Transmission Dynamics

- Transmission Models for SARS-COV virus
- Superspreading and the effect of individual variation on disease emergence



# Results: Popular Softwares and ML Models

## Popular techniques and algorithms used for SARS studies





# Results: Popular Softwares and Packages

## Popular Softwares and Packages:

- Only one of the papers mentions a coding language, which is Java and C. All of the other papers don't mention how the models were created (i.e. using a specific program or a programming language) and no links were found to any type of code/repositories. We believe that's due to the time period of the virus as before mentioned. Majority of the papers that were found date back more than 15 years.

## Quality of papers and dataset:

- Large part of the papers don't include their datasets. The ones that do include links to datasets, the links are out of date and aren't reachable.
- External Datasets were found to fill the gap of data. The data that has been found is 1) Case/Deaths Time Histories for many different parts of the world and 2) X-Ray images from SARS patients.

# Results: Datasets from Machine Learning Publications

	Dataset	Link
1	Table with global data for the SARS epidemic	<a href="https://github.com/WL-Biol185-ShinyProjects/sars-project/blob/master/SARS%20data.xlsx">https://github.com/WL-Biol185-ShinyProjects/sars-project/blob/master/SARS%20data.xlsx</a>
2	The file contains day by day no. from March to July 2003 across the world.	<a href="https://www.kaggle.com/imdevskp/comments">https://www.kaggle.com/imdevskp/comments</a>
3	Final summary data from across the world	<a href="https://www.kaggle.com/imdevskp/comments">https://www.kaggle.com/imdevskp/comments</a>
4	Chest X Ray dataset for SARS	<a href="https://github.com/mlmed/torchxrayvision">https://github.com/mlmed/torchxrayvision</a>

# COVID-19

---

Mengyao He, Om Vaghasia



# Introduction

## Literature Review:

- Researchers are using Machine Learning techniques to apply to problems related to COVID-19, from ***diagnosing early symptoms*** using Chest X-rays to internet search queries to ***estimate the effect on mental health of social distancing***.
- Most of the research is focused on classification and prediction.
- The techniques that have been mostly used are :
  - K-means, K-nearest neighbors
  - Random Forest and tree-based models in general
  - Neural Networks and Deep Learning
- They mostly use ***time-series data*** of ***infections/deaths/recoveries, chest X-rays and CT scan***, that were collected in different cities around the world.
- Heavily subsidized research topic, leading to a vast amount of research papers and datasets.

# Methods to Conduct Literature Reviews

## PubMed Keywords

- Pubmed search engine:
  - <https://pubmed.ncbi.nlm.nih.gov/>
- COVID-19:
  - COVID19 Machine Learning
  - SARS-CoV-2 Machine Learning
- Results:
  - 59 publications
  - 34 publications datasets

## Other Methods

- [Harvard Dataverse](#)
- [Kaggle](#)
- Github repositories
- Links provided within other sections of the papers

# Results: Major Groups in ML-based Covid-19 Studies

## Screening

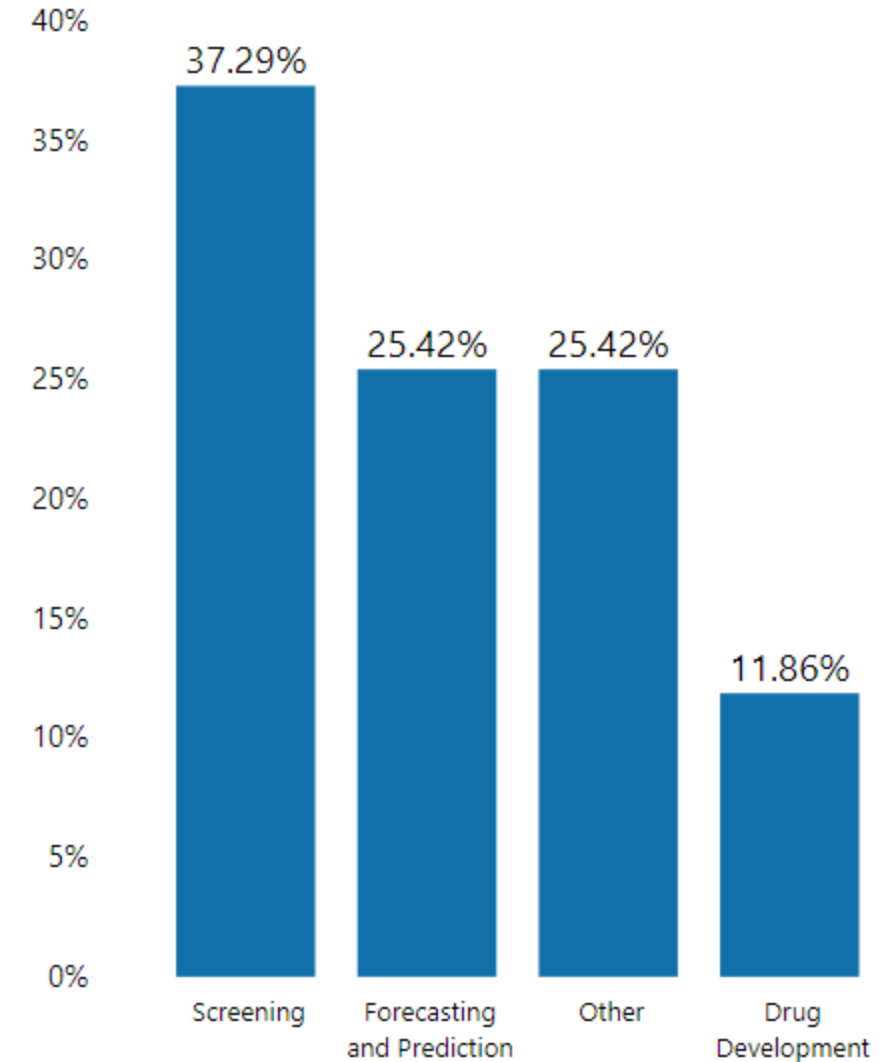
- Extract radiological features for timely and accurate COVID-19 diagnosis from CT images and X rays.
- Distinguish COVID-19 from community acquired Pneumonia and other lung diseases.

## Forecasting & Prediction

- Forecast the trajectory of the outbreak and quantify the risk of death.
- Show feasibility and accuracy for predicting hospital stay in COVID-19 patients
- Protein structure predictions

## Drug Development

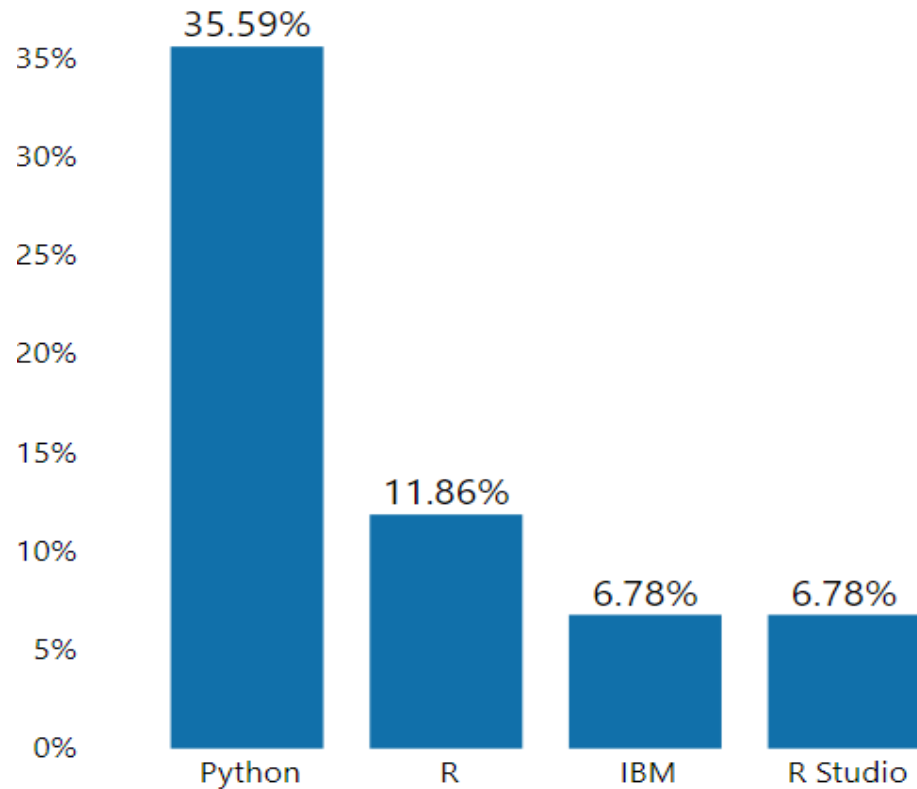
- Computational biology and medicines perspective
- Generate novel drug compounds



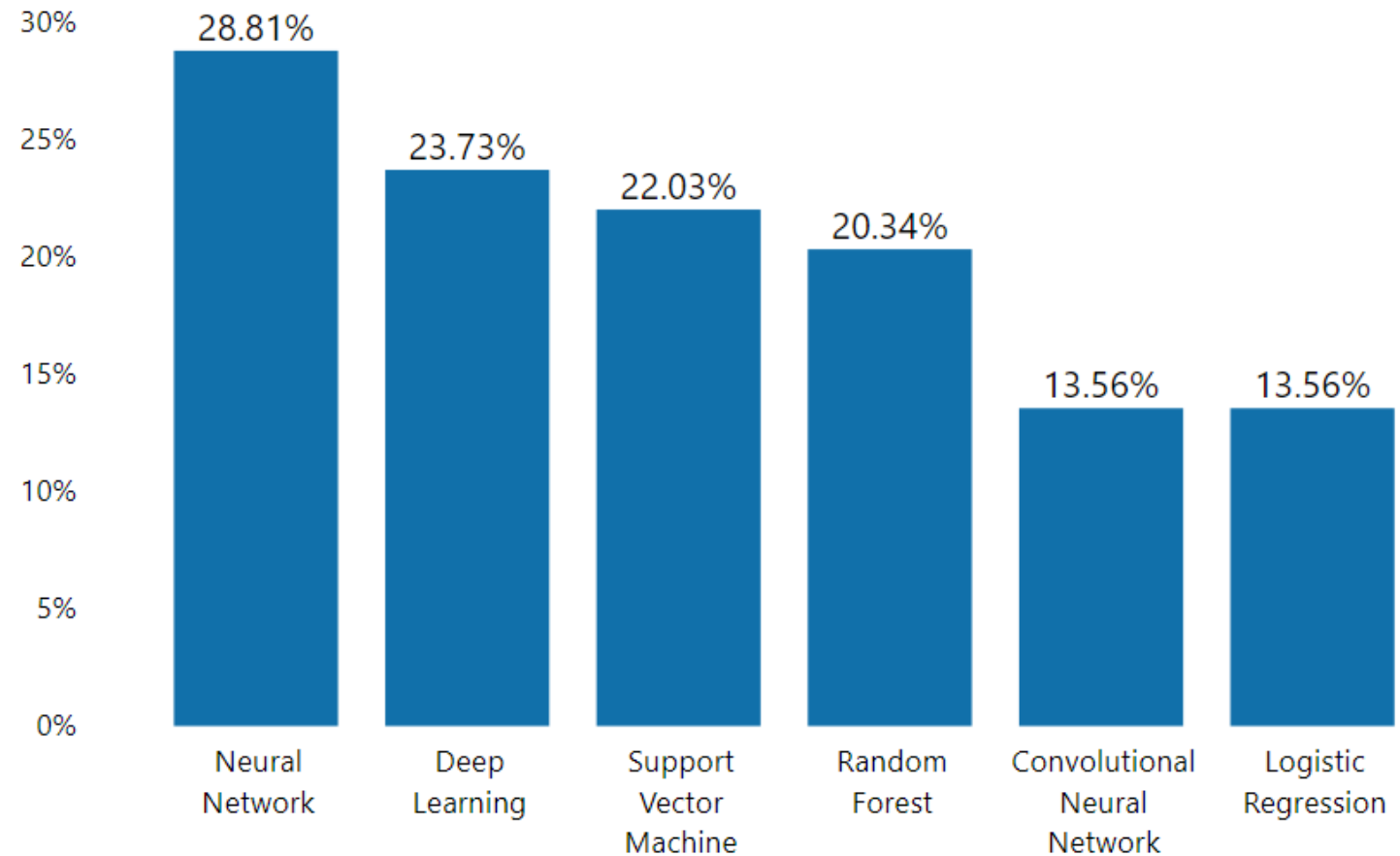


# Results: Popular Softwares and ML Models

**Popular softwares used for ML-based studies for COVID-19**



**Popular ML models and algorithms used for COVID-19 studies**



# Results: Datasets from COVID-19 Machine Learning Publications

	Dataset	Link		Dataset	Link
1	Pulmonary Chest CT Scans	<a href="https://github.com/bkong999/COVNet">https://github.com/bkong999/COVNet</a>	10	Baidu migration	<a href="http://qianxi.baidu.com/">http://qianxi.baidu.com/</a>
2	COVID-19 Affected Cases	<a href="https://www.kaggle.com/sudalairajkumar/novel-corona-virus-2019-dataset">https://www.kaggle.com/sudalairajkumar/novel-corona-virus-2019-dataset</a>	11	Real-time cases of COVID-19	<a href="https://news.qq.com/zt2020/page/feiyan.htm?ADTAG=area">https://news.qq.com/zt2020/page/feiyan.htm?ADTAG=area</a>
3	COVID-Chestxray-Dataset	<a href="https://github.com/ieee8023/covid-chestxray-dataset">https://github.com/ieee8023/covid-chestxray-dataset</a>	12	Situation report 2020	<a href="http://www.nhc.gov.cn/xcs/yqtb/list_gzbd.shtml">http://www.nhc.gov.cn/xcs/yqtb/list_gzbd.shtml</a>
4	CSSEGISandData	<a href="https://github.com/CSSEGISandData/COVID-19">https://github.com/CSSEGISandData/COVID-19</a>	13	Combatting SARS 2003	<a href="http://news.sohu.com/57/26/subject206252657.shtml">http://news.sohu.com/57/26/subject206252657.shtml</a>
5	Protein Sequences	<a href="https://bigd.big.ac.cn/ncov">https://bigd.big.ac.cn/ncov</a>	14	Dataset of SARS-CoV-2 Genome	<a href="https://data.mendeley.com/datasets/nvk5bf3m2f/1">https://data.mendeley.com/datasets/nvk5bf3m2f/1</a>
6	Proteome Data	<a href="https://www.iprox.org/page/ProjectFileList.html?projectId=IPX0002106000">https://www.iprox.org/page/ProjectFileList.html?projectId=IPX0002106000</a>	15	COVID-Classfier	<a href="https://github.com/abzargar/COVID-Classfier/tree/master/dataset">https://github.com/abzargar/COVID-Classfier/tree/master/dataset</a>
7	MLDSP-GUI	<a href="https://sourceforge.net/projects/mldsp-gui/files/COVID19Dataset/">https://sourceforge.net/projects/mldsp-gui/files/COVID19Dataset/</a>	16	NYT Data Collection	<a href="https://github.com/mihirpsu/covid_19/tree/master/Data_Collection">https://github.com/mihirpsu/covid_19/tree/master/Data_Collection</a>
8	COVID-19 Cases Dataset	<a href="https://www.kaggle.com/imdevskp/corona-virus-report">https://www.kaggle.com/imdevskp/corona-virus-report</a>	17	COVID-19 Brodinlab	<a href="https://kl.app.box.com/s/sby0jesyu23a65cbgv51vpbzqjdmjpr1">https://kl.app.box.com/s/sby0jesyu23a65cbgv51vpbzqjdmjpr1</a>
9	COVID-19 Global Wealth	<a href="https://www.kaggle.com/winterpierre91/covid19-global-weather-data">https://www.kaggle.com/winterpierre91/covid19-global-weather-data</a>	18	UCSDAI4H COVID-19 CT	<a href="https://github.com/UCSD-AI4H/COVID-CT">https://github.com/UCSD-AI4H/COVID-CT</a>

# Results: Datasets from COVID-19 Machine Learning Publications

	Dataset	Link		Dataset	Link
19	Confirmed COVID-19 cases	<a href="https://doi.org/10.6084/m9.figshare.12030363.v1">https://doi.org/10.6084/m9.figshare.12030363.v1</a>	27	GCCR001	<a href="https://osf.io/a3vkw/">https://osf.io/a3vkw/</a>
20	Multimedia component	<a href="https://www.ncbi.nlm.nih.gov/pmc/articles/PMC7186211/bin/mmc1.csv">https://www.ncbi.nlm.nih.gov/pmc/articles/PMC7186211/bin/mmc1.csv</a>	28	100 axial CT images	<a href="http://www.salhospital.com/">http://www.salhospital.com/</a>
21	Health Report from China CDC	<a href="https://github.com/midas-network/COVID-19/tree/master/data/cases/china">https://github.com/midas-network/COVID-19/tree/master/data/cases/china</a>	29	Public Chest X-Ray (CXR) datasets	<a href="https://www.isi.uu.nl/Research/Databases/SCR/download.php">https://www.isi.uu.nl/Research/Databases/SCR/download.php</a>
22	CoronaHack CXR Dataset	<a href="https://www.kaggle.com/praveengovi/coronahack-chest-xraydataset">https://www.kaggle.com/praveengovi/coronahack-chest-xraydataset</a>	30	Montgomery County X-ray Set	<a href="https://academictorrents.com/details/ac786f74878a5775c81d490b23842fd4736bfe33">https://academictorrents.com/details/ac786f74878a5775c81d490b23842fd4736bfe33</a>
23	CXR Metadata	<a href="https://github.com/jongcye/Deep-Learning-COVID-19-on-CXR-using-Limited-Training-Data-Sets/blob/master/metadata.xls">https://github.com/jongcye/Deep-Learning-COVID-19-on-CXR-using-Limited-Training-Data-Sets/blob/master/metadata.xls</a>	31	JSRT Database	<a href="http://db.jsrt.or.jp/eng.php">http://db.jsrt.or.jp/eng.php</a>
24	Confirmed COVID-19 cases	<a href="https://doi.org/10.6084/m9.figshare.12030363.v1">https://doi.org/10.6084/m9.figshare.12030363.v1</a>	32	GCCR001	<a href="https://osf.io/a3vkw/">https://osf.io/a3vkw/</a>
25	Multimedia component	<a href="https://www.ncbi.nlm.nih.gov/pmc/articles/PMC7186211/bin/mmc1.csv">https://www.ncbi.nlm.nih.gov/pmc/articles/PMC7186211/bin/mmc1.csv</a>	33	100 axial CT images	<a href="http://www.salhospital.com/">http://www.salhospital.com/</a>
26	COVID Chestray Dataset	<a href="https://github.com/leee8023/covid-chestxray-dataset">https://github.com/leee8023/covid-chestxray-dataset</a>	34	Serial chest radiographs graphs	<a href="https://www.ncbi.nlm.nih.gov/pmc/articles/PMC7135076/figure/fig1/">https://www.ncbi.nlm.nih.gov/pmc/articles/PMC7135076/figure/fig1/</a>

# SARS (2003) and COVID-19 Conclusions

## SARS (2003)

- Most of the methodologies that were used in the literature are outdated. The most common techniques found are Neural Networks, Decision Trees and Stochastic Models.
- In terms of data all of the articles weren't useful in providing data. External Datasets were found.
- Majority of papers date back more than 15 years, opportunity for researchers to use newer ML techniques on the SARS epidemic.

## COVID-19

- The most common techniques found are Neural Networks, Deep Learning, and SVM.
- K-Nearest Neighbors is observed to be a preferred method for clustering.
- Size of datasets for studies on Screening seems to be an issue, particularly for Deep Learning and Convolutional Neural Network algorithms. Transfer Learning can be explored as a potential solution.



# Future Work

## SARS (2003)

Researchers could use the SARS datasets that have been found and the methods from the COVID-19 Machine Learning papers to implement more up-to-date techniques and see how they compare to the ones that were used 15 years ago.

## COVID-19

Publications on ML-based COVID-19 studies are being conducted and uploaded continuously which can be studied and added to the current database.

Trends regarding the type of studies being conducted - Screening, Forecasting and Prediction, and Drug Development - can be tracked to identify some of the needs.

# References

- Lalmuanawma S, Hussain J, Chhakchhuak L. Applications of machine learning and artificial intelligence for Covid-19 (SARS-CoV-2) pandemic: A review. *Chaos Solitons Fractals*. 2020;139:110059. doi:10.1016/j.chaos.2020.110059
- Kumar A, Gupta PK, Srivastava A. A review of modern technologies for tackling COVID-19 pandemic. *Diabetes Metab Syndr*. 2020;14(4):569-573. doi:10.1016/j.dsx.2020.05.008

# Thank You!

*Thank you for this opportunity*, especially to Ravi and Naomi. We learned about FNL and about collaborating to solve real machine learning problems in the medical science. We've learned a lot and we hope to be able to collaborate with FNL again.

We have identified *70 publications and 38 publicly available datasets* for machine learning studies in the SARS space.

*Any questions?*