



# 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



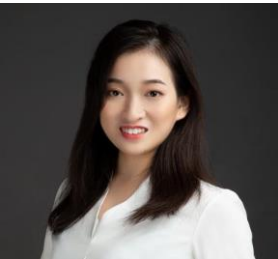
Jim Hu



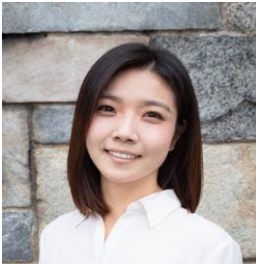
Niranjana Moleyar



Malin Ortenblad



Kerry Hu



Jie Chen



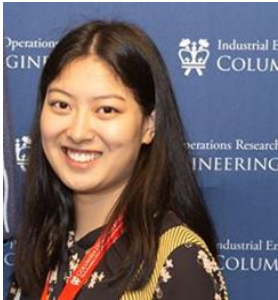
Mengyao He



Om Vaghasia



Panagiotis Misirlis



Jiaxi Zhou



Xinyao Wang



Qinwei Zhang



Yue Hu



Zihui Zhou



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 Top-10 Commonly Occurring Cancer Datasets for ML

Team Members: Jim Hu, Kerry Hu, Mahitha Kotipalli, Malin Ortenblad, Niranjana Moleyar  
Affiliation: IEOR Department, Columbia University  
Sep. 8, 2020

# Table of content

- **Project Overview**
- Methodology and Summary
- Results
  - Skin and Breast Cancer
  - Kidney and Colon Cancer
  - Spinal Cord, Cranial Nerves and Ovary
  - Bronchus, Lung and Prostate Gland
- Next Steps

# Project Team



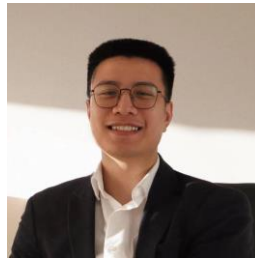
Mahitha Kotipalli  
Columbia University, MSMSE  
Project 1: Bronchus, Lung and Prostate Cancer  
Prior experience in healthcare consulting. Worked as Business Analytics MBA intern with Bayer, Business Insights team this summer



Malin Ortenblad  
Columbia University, MSBA  
Project 1: Skin and Breast Cancer  
Previous work experience in healthcare consulting, spent the summer at Bayer in their Business Insights department as a Business Analytics MBA.



Kerry Hu  
Columbia University, MSOR  
Project 1: Spinal Cord, Cranial Nerves and Ovary  
Prior experience in financial services. Interned at Credit Suisse as quantitative researcher this summer



Jim Hu  
Columbia University, MSBA  
Project 1: Bronchus, Lung and Prostate Cancer  
Previous Data Science Researcher at Point72, Business Analyst at Bayer Women's Health, Business Insights



Niranjana Moleyar  
Columbia University, MSMSE  
Kidney and Colon Cancer  
Prior experience in Recommendation Systems and Convolutional Neural Networks for image comparison. Worked with Rent the Runway to optimise their inventory



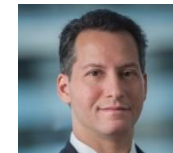
Jie Chen  
Columbia University, MSOR  
Project 1: Hematopoietic and reticuloendothelial (bone marrow) and Pancreas  
Prior experience in financial service. Interned at Bank of America as GBAM associate intern.



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



# Project Overview: What Human Cancer Datasets (Biomolecular/Drug/Phenotype) are Available for Machine-Learning?

## Objective

- Using PubMed to identify available datasets for machine learning in the oncology space, focusing on;
  - Breast, Skin, Lung, Bronchus, Prostate, Kidney, Colon, and Spinal Cord Cancer

## Process

- Published work summarized in excel trackers with links to publication, abstract, title, dataset, and methodology
- Results available at the following GitHub repositories:
  - Skin and breast cancer: [https://github.com/Ortenblad/ColUniv-FNL-BIDS-Project1\\_BreastandSkinCancer](https://github.com/Ortenblad/ColUniv-FNL-BIDS-Project1_BreastandSkinCancer)
  - Bronchus, Lung and Prostate cancer: [https://github.com/Jim-Hu/ColUniv-FNL-BIDS-Project1\\_LungBronchusProstate](https://github.com/Jim-Hu/ColUniv-FNL-BIDS-Project1_LungBronchusProstate)
  - Kidney and Colon Cancer: <https://github.com/NiranjanaMoleyar/FNL-BIDS-Project1-KidneyAndColonCancer>
  - Spinal Cord, Cranial Nerves and Ovary: [https://github.com/foevee/FNL-BIDS-Project1-SpinalCord\\_CranialNerves\\_Ovary](https://github.com/foevee/FNL-BIDS-Project1-SpinalCord_CranialNerves_Ovary)

# Table of content

- **Project Overview**
- **Methodology and Summary**
- **Results**
  - Skin and Breast Cancer
  - Kidney and Colon Cancer
  - Spinal Cord, Cranial Nerves and Ovary
  - Bronchus, Lung and Prostate Gland
- **Next Steps**



# In addition to searching for papers, many authors were contacted and bibliographies were used to identify additional research and data

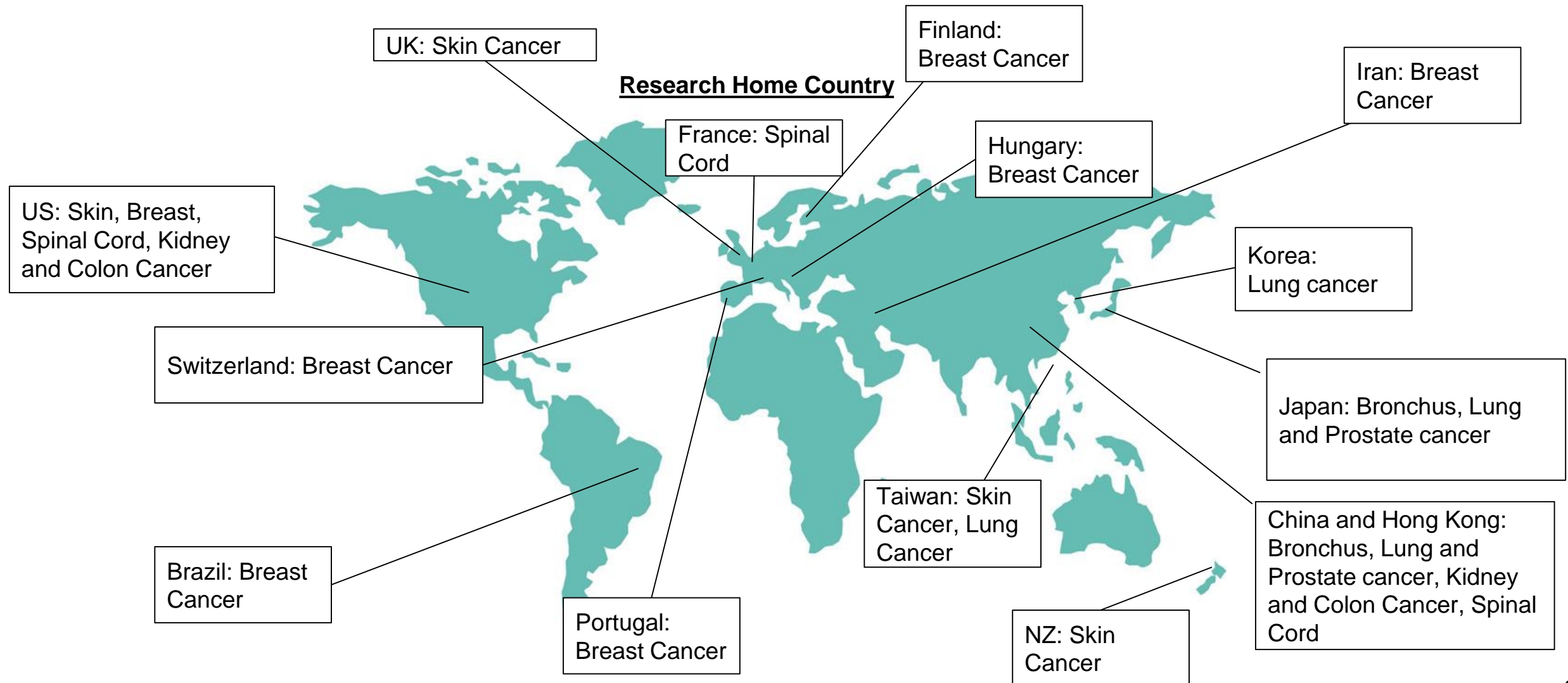
## Search Terms

- Pubmed search engine link: <https://pubmed.ncbi.nlm.nih.gov/>
- Skin and Breast:
  - Machine Learning Skin Cancer
  - Machine Learning Breast Cancer
- Bronchus, Lung and Prostate Cancer:
  - Machine Learning Bronchus Cancer
  - Machine Learning / Deep Learning Lung Cancer
  - Machine Learning Prostate Cancer
- Kidney and Colon Cancer:
  - Kidney Cancer Machine Learning / Artificial Intelligence
  - Colon Cancer Machine Learning / Artificial Intelligence
- Spinal Cord, Cranial Nerves and Ovary:
  - Spinal Cord and Machine Learning / Neural Networks
  - Cranial Nerves and Machine Learning / Neural Networks
  - Ovary and Machine Learning / Neural Networks

## Additional Steps

- Identify additional publications from footnotes and bibliography of papers
- Emailed authors requesting dataset (response rate ~5%)
- Looked up the identified papers at the Plos One journal since they provide data sources frequently

# Based on the papers we gathered, research is conducted globally



# Issues with gathering datasets leveraged in publications

## Common Data Sources

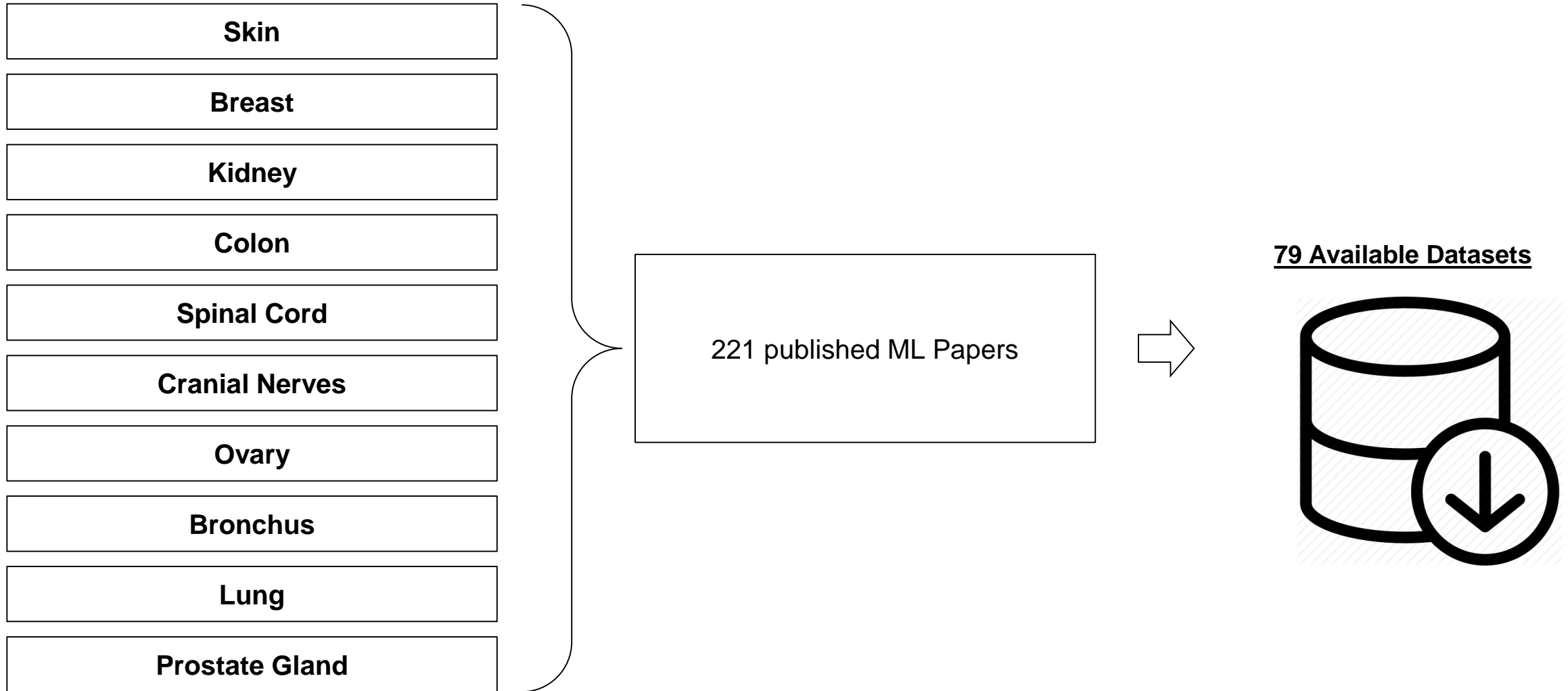
- Most of the research is used typically comes from the following sources
  - a. Publicly available datasets (e.g. on the website of **National Cancer Institute**)
  - b. Datasets published on the website of different **universities' cancer research groups**
  - c. Datasets available from competitions
  - d. Data gathered by the researcher for the specific paper (hardest to acquire)

## Issues faced to access the data

- Data owned by universities and other research groups is often difficult to access because
  - a. **Regulations** in the country where research is carried out
  - b. Requirement of **contract** between the parties to share data
  - c. Some data has to be **deleted** after five years (limiting data availability to research published in last five years)
- Low author response rates



# Project team has summarized 221 publications, and identified 79 datasets across 10 cancer types



# Table of content

- Project Overview
- Methodology
- **Results**
  - **Skin and Breast Cancer**
  - Kidney and Colon Cancer
  - Spinal Cord, Cranial Nerves and Ovary
  - Bronchus, Lung and Prostate Gland
- Next Steps

# For breast and skin cancer the most cost common use of ML was for classification and diagnosis purposes

## Breast Cancer

### Use cases and objectives for ML research:

- Classification (most common):
  - Create machine to diagnose breast cancer
  - Outlining machine learning algorithms in use to detect breast cancer in mammogram images
  - Automate the method to detect cancer using ultrasounds
  - Comparison of different algorithms in diagnosing breast cancer
  - Improving diagnosis accuracy
- Patient Journey:
  - Predicting regression patterns in breast cancer patients
  - Predicting Tumor growth
  - Symptom analysis
- Sample: 50

## Skin Cancer

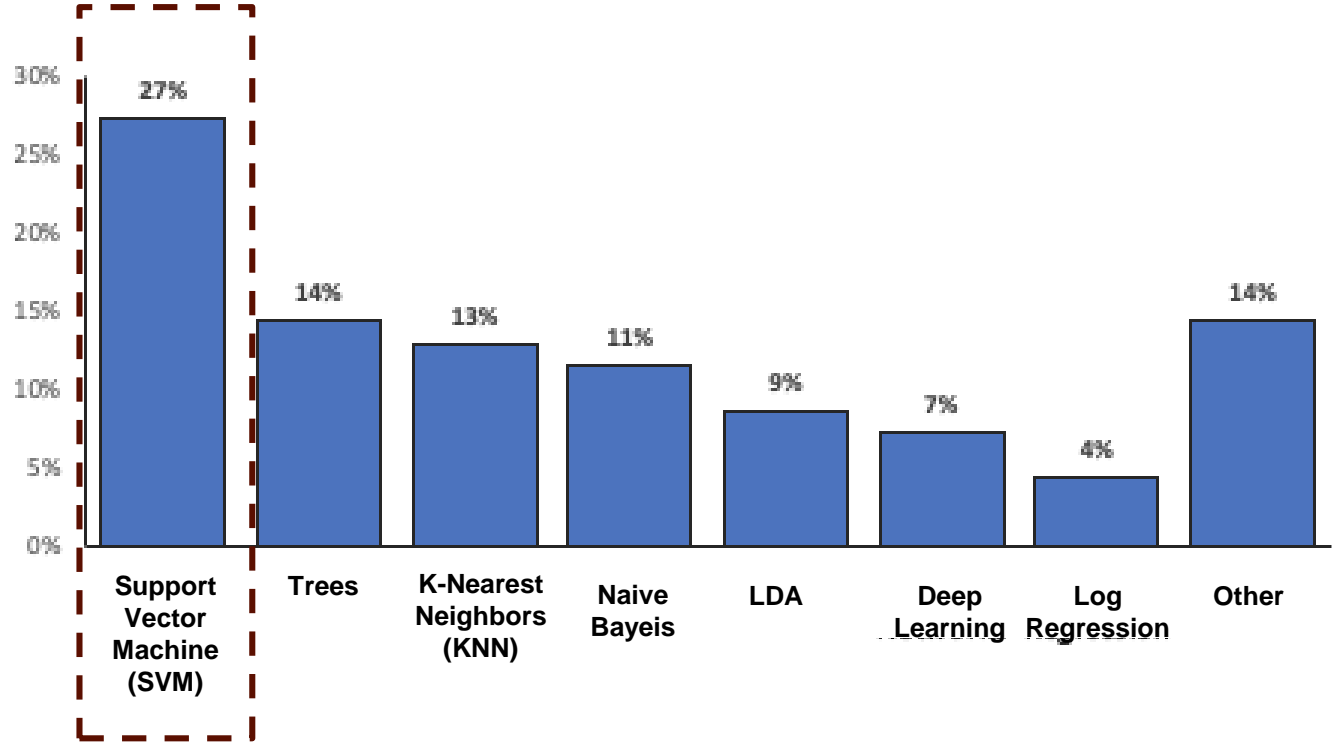
### Use cases and objectives for ML research:

- Classification:
  - Comparison of AI and dermatologist in identifying skin cancer
  - Development and testing of 3D image scanning
  - Evaluation of image quality on diagnosis accuracy
  - Skin cancer diagnosis: ML, dermatologists, vision enhancing tools
  - Clinical feature selection and discovery of new biomarkers
  - Differentiative skin cancer disease classification
  - Smartphone app classifier accuracy
- Patient Journey:
  - Patient pathway utilization: Survivorship care plans
- Sample:39

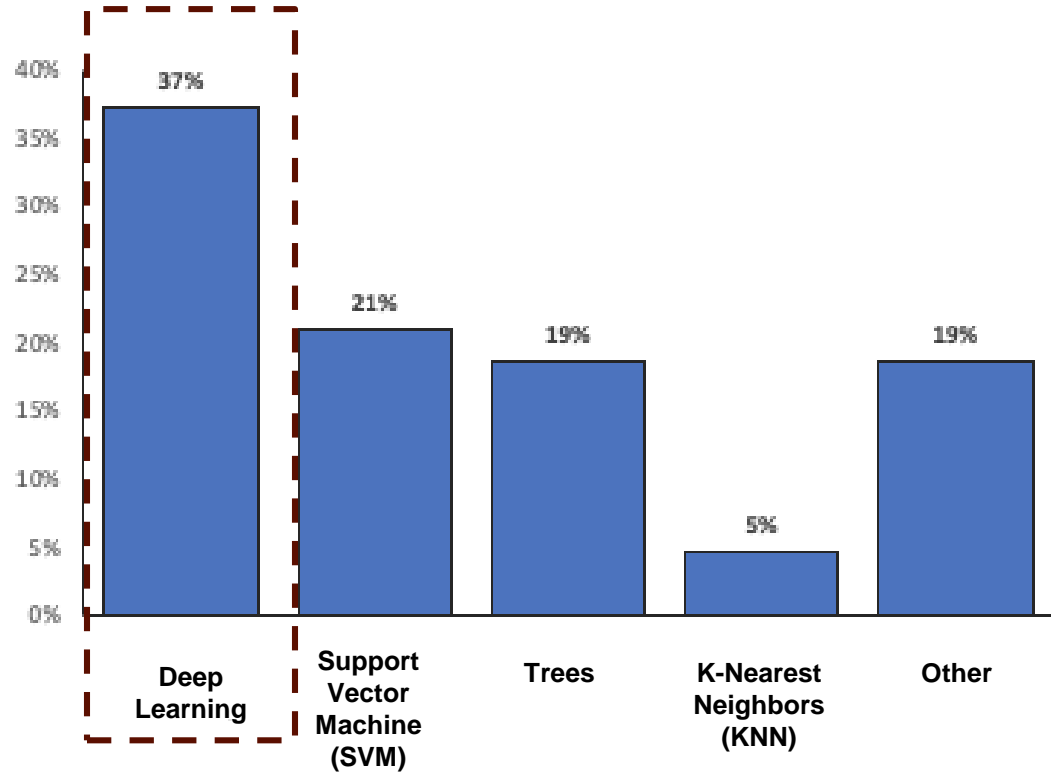


# The most common algorithms to use in breast and skin cancer research are SVM and Neural Networks, respectively

Algorithms Used for ML in Breast Cancer



Algorithms Used for ML in Skin Cancer



# 24 Datasets were identified from Breast Cancer Machine learning Publications

	Dataset	Link		Dataset	Link
1	GSE59198	<a href="https://www.omicsdi.org/dataset/geo/GSE59198">https://www.omicsdi.org/dataset/geo/GSE59198</a>	13	Mouse-Mammary	<a href="https://wiki.cancerimagingarchive.net/display/Public/Mouse-Mammary">https://wiki.cancerimagingarchive.net/display/Public/Mouse-Mammary</a>
2	Research Patient Data Registry	<a href="https://rc.partners.org/about/who-we-are-risc/research-patient-data-registry">https://rc.partners.org/about/who-we-are-risc/research-patient-data-registry</a>	14	TCGA-BRCA	<a href="https://portal.gdc.cancer.gov/projects/TCGA-BRCA">https://portal.gdc.cancer.gov/projects/TCGA-BRCA</a>
3	MIAS	<a href="https://www.kaggle.com/kmader/mias-mammography">https://www.kaggle.com/kmader/mias-mammography</a>	15	QIN Breast DCE-MRI	<a href="https://wiki.cancerimagingarchive.net/display/Public/QIN+Breast+DCE-MRI">https://wiki.cancerimagingarchive.net/display/Public/QIN+Breast+DCE-MRI</a>
4	DDSM	<a href="https://www.kaggle.com/skooch/ddsm-mammography">https://www.kaggle.com/skooch/ddsm-mammography</a>	16	BREAST-DIAGNOSIS	NA
5	CBIS-DDSm	<a href="https://wiki.cancerimagingarchive.net/display/Public/CBIS-DDSM">https://wiki.cancerimagingarchive.net/display/Public/CBIS-DDSM</a>	17	RIDER Breast MRI	<a href="https://wiki.cancerimagingarchive.net/display/Public/RIDER+Breast+MRI">https://wiki.cancerimagingarchive.net/display/Public/RIDER+Breast+MRI</a>
6	ISPY1	<a href="https://wiki.cancerimagingarchive.net/display/Public/ISPY1">https://wiki.cancerimagingarchive.net/display/Public/ISPY1</a>	18	BCDR	<a href="https://bcdr.eu/information/about">https://bcdr.eu/information/about</a>
7	Breast-MRI-NACT-Pilot	<a href="https://wiki.cancerimagingarchive.net/display/Public/Breast-MRI-NACT-Pilot">https://wiki.cancerimagingarchive.net/display/Public/Breast-MRI-NACT-Pilot</a>	19	TCGA-BRCA	<a href="https://portal.gdc.cancer.gov/projects/TCGA-BRCA">https://portal.gdc.cancer.gov/projects/TCGA-BRCA</a>
8	QIN-Breast	<a href="https://wiki.cancerimagingarchive.net/display/Public/QIN-Breast">https://wiki.cancerimagingarchive.net/display/Public/QIN-Breast</a>	20	BreakHis	<a href="https://web.inf.ufpr.br/vri/databases/breast-cancer-histopathological-database-breakhis/">https://web.inf.ufpr.br/vri/databases/breast-cancer-histopathological-database-breakhis/</a>
9	Wisconsin Breast Cancer (original) datasets	<a href="https://archive.ics.uci.edu/ml/datasets/breast+cancer+wiscconsin+(original)">https://archive.ics.uci.edu/ml/datasets/breast+cancer+wiscconsin+(original)</a>	21	Inbreast	<a href="https://github.com/wentaozhu/deep-mil-for-whole-mammogram-classification/issues/12">https://github.com/wentaozhu/deep-mil-for-whole-mammogram-classification/issues/12</a>
10	ICAR 2018	<a href="https://iciar2018-challenge.grand-challenge.org/Dataset/">https://iciar2018-challenge.grand-challenge.org/Dataset/</a>	22	Finprog Cancer Dataset	<a href="http://finprog.org/data_set2.asp">finprog.org/data_set2.asp</a>
11	GWAS	<a href="https://cbi.oml.gov/data">https://cbi.oml.gov/data</a>	23	Metabolomics data of the 162 metabolites	<a href="https://www.sciencedirect.com.ezproxy.cul.columbia.edu/science/article/pii/S1874391913005113?via%3Dihub">https://www.sciencedirect.com.ezproxy.cul.columbia.edu/science/article/pii/S1874391913005113?via%3Dihub</a>
12	271 breast cancer samples	<a href="https://www.ncbi.nlm.nih.gov.ezproxy.cul.columbia.edu/geo/">https://www.ncbi.nlm.nih.gov.ezproxy.cul.columbia.edu/geo/</a>	24	Mini-MIAS	<a href="http://peipa.essex.ac.uk/info/mias.html">http://peipa.essex.ac.uk/info/mias.html</a>

## 20 Datasets were identified from Skin Cancer Machine learning Publications

	Dataset	Link		Dataset	Link
1	Innternational Skin Image Collaboration	<a href="https://www.isic-archive.com/#!/topWithHeader/wideContentTop/main">https://www.isic-archive.com/#!/topWithHeader/wideContentTop/main</a>	11	COSMIC	<a href="https://cancer.sanger.ac.uk/cosmic">https://cancer.sanger.ac.uk/cosmic</a>
2	NHIRD	<a href="https://nhird.nhri.org.tw/en/">https://nhird.nhri.org.tw/en/</a>	12	DermNEt	<a href="https://dermnetnz.org/">https://dermnetnz.org/</a>
3	SCP Survey Data	<a href="http://www.scpwiki.com/2018-survey-results">http://www.scpwiki.com/2018-survey-results</a>	13	ImageNet	<a href="http://www.image-net.org/">http://www.image-net.org/</a>
4	digital dermoscopic database	<a href="https://www.sec.gov/Archives/edgar/data/1051514/000095012305009834/y09078a3sv1za.htm">https://www.sec.gov/Archives/edgar/data/1051514/000095012305009834/y09078a3sv1za.htm</a>	14	TCGA-SKCM	<a href="https://portal.gdc.cancer.gov/projects/TCGA-SKCM">https://portal.gdc.cancer.gov/projects/TCGA-SKCM</a>
5	ISIC Challenge dataset	<a href="https://challenge2020.isic-archive.com/">https://challenge2020.isic-archive.com/</a>	15	Asan dataset	<a href="https://figshare.com/articles/Asan_and_Hallym_Dataset_Thumbnails/_5406136">https://figshare.com/articles/Asan_and_Hallym_Dataset_Thumbnails/_5406136</a>
6	MED-NODE	<a href="http://www.cs.rug.nl/~imaging/databases/melanoma_naevi/">http://www.cs.rug.nl/~imaging/databases/melanoma_naevi/</a>	16	Edinburgh dataset	<a href="https://licensing.eri.ed.ac.uk/i/software/dermofit-image-library.html">https://licensing.eri.ed.ac.uk/i/software/dermofit-image-library.html</a>
7	Atlas Derm dataset	<a href="https://www.dermatlas.net/">https://www.dermatlas.net/</a>	17	Hallym Dataset	<a href="https://figshare.com/articles/Asan_and_Hallym_Dataset_Thumbnails/_5406136">https://figshare.com/articles/Asan_and_Hallym_Dataset_Thumbnails/_5406136</a>
8	DermIS	<a href="https://www.dermis.net/dermisroot/en/home/index.htm">https://www.dermis.net/dermisroot/en/home/index.htm</a>	18	ILSVRC; ISBI 2016 (Gutman	<a href="#">ISBI 2016 (Gutmanet al., 2016).</a>
9	Derm101 (No longer available under this name)	<a href="https://www.emailmeform.com/builder/form/Ne0j8da9bb7U4h6t1f">https://www.emailmeform.com/builder/form/Ne0j8da9bb7U4h6t1f</a>	19	HAM1000	<a href="https://dataverse.harvard.edu/dataset.xhtml?persistentId=doi:10.7910/DVN/DBW86T">https://dataverse.harvard.edu/dataset.xhtml?persistentId=doi:10.7910/DVN/DBW86T</a>
10	GSE122703	<a href="https://www.ncbi.nlm.nih.gov/geo/query/acc.cgi?acc=GSE122703">https://www.ncbi.nlm.nih.gov/geo/query/acc.cgi?acc=GSE122703</a>	20	992 dermoscopy images	<a href="http://www.bccancer.bc.ca/HPI/SkinCancerAtlas/Melanoma/default.htm">http://www.bccancer.bc.ca/HPI/SkinCancerAtlas/Melanoma/default.htm</a> <a href="http://www.cancer.org/cancer/skin-cancer/melanoma/detailedguidelines/melanoma-skin-cancer-key-statistics">http://www.cancer.org/cancer/skin-cancer/melanoma/detailedguidelines/melanoma-skin-cancer-key-statistics</a> <a href="http://www.dermnet.org.nz/files/online/melanoma-aim-el-is.html">http://www.dermnet.org.nz/files/online/melanoma-aim-el-is.html</a> <a href="http://www.medicare.nsw.edu.au/en/MeDicof/medicof/dermatology/medicof-environment-1.htm">http://www.medicare.nsw.edu.au/en/MeDicof/medicof/dermatology/medicof-environment-1.htm</a>



# Table of content

- Project Overview
- Methodology
- **Results**
  - Skin and Breast Cancer
  - **Kidney and Colon Cancer**
  - Spinal Cord, Cranial Nerves and Ovary
  - Bronchus, Lung and Prostate Gland
- Next Steps

# Most common use of ML in Kidney and Colon cancer is in identification, prognosis and progress prediction

## Kidney Cancer

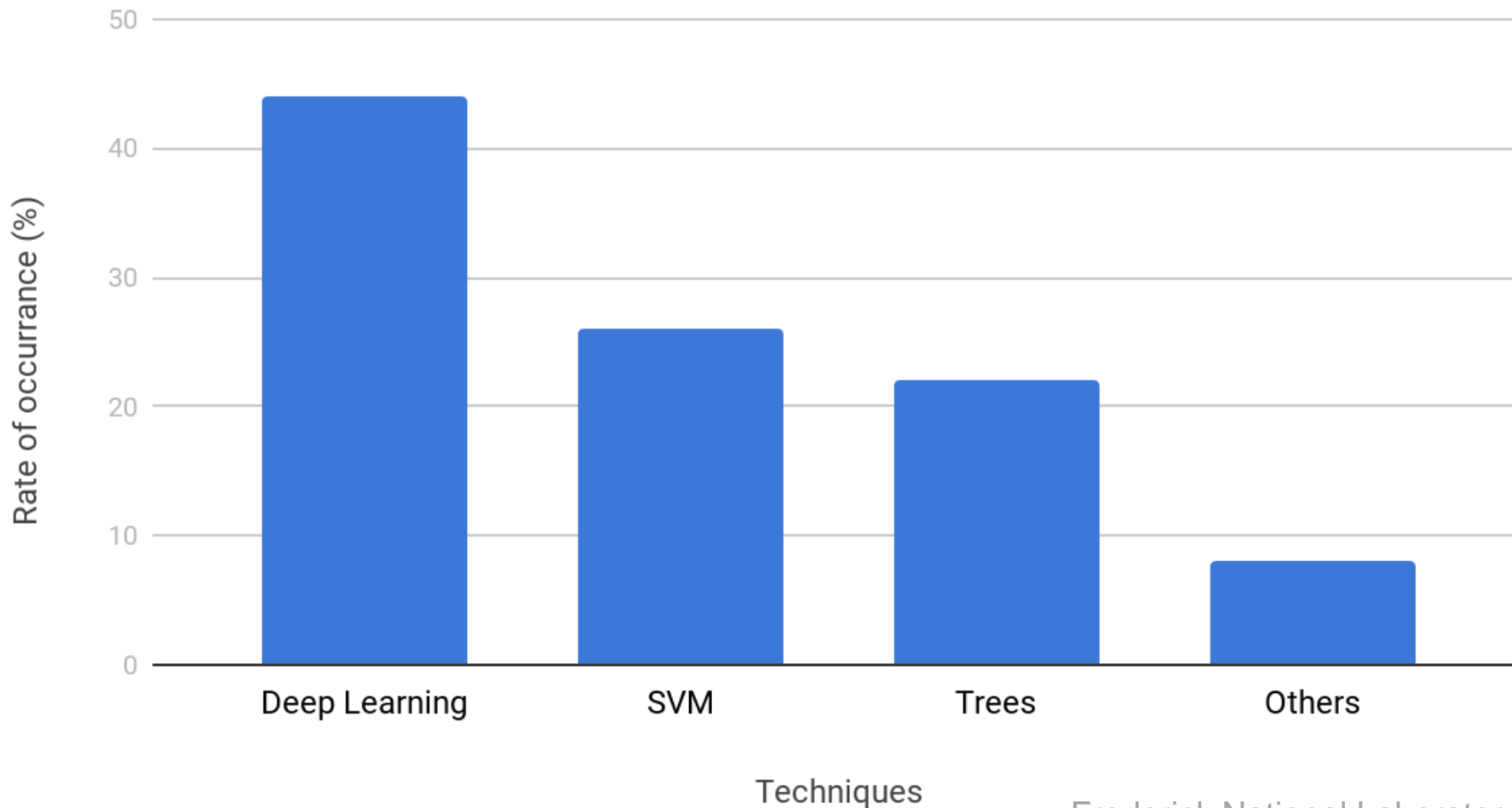
- ML based CT analysis
- Prognosis and survival prediction
- Automated detection
- Radiation dose estimation
- Prediction of the nuclear grade of renal cell carcinoma

## Colon Cancer

- Prediction of progression of colon cancer and survival
- Expanding TNM for cancers
- Establishing gene expression signature for prognosis
- Automated detection
- Classification of cancer types

# Deep Learning is the most widely used technique followed by SVM

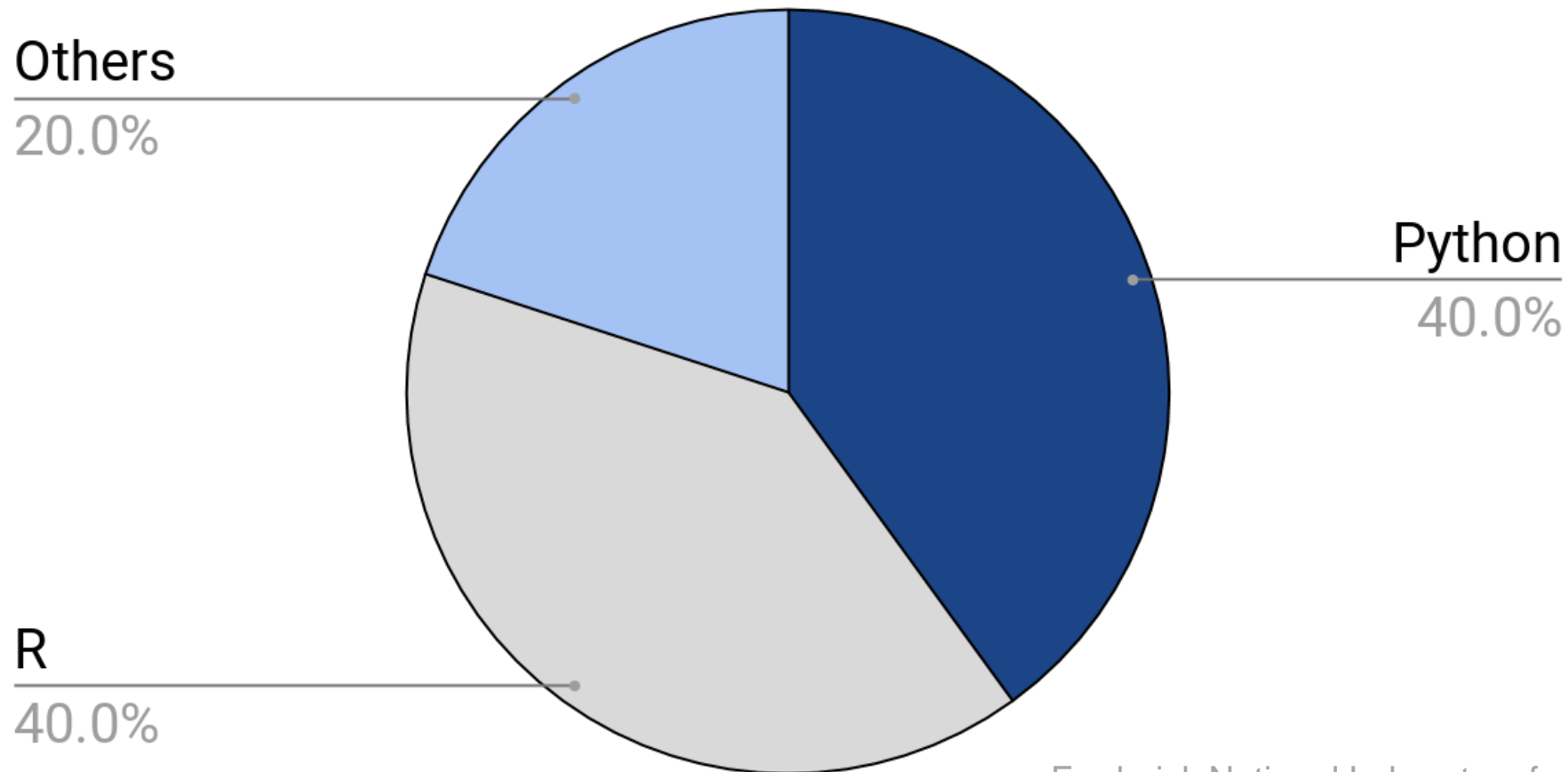
Techniques used for ML in Kidney and Colon Cancer Research



# Python and R are predominantly used in research papers for Kidney and Colon cancer

Python is mainly used for deep learning, whereas R is predominantly used for more traditional ML techniques.

Softwares used for data analysis in Kidney and Colon Cancer research



# List of common public data repositories and datasets

Dataset	Cancer Type	Link
TCGA (The Cancer Genome Atlas) data repository	Any	<a href="https://cancergenome.nih.gov/">https://cancergenome.nih.gov/</a>
Harvard Dataverse	Any	<a href="https://dataverse.harvard.edu/">https://dataverse.harvard.edu/</a>
Gene Expression Omnibus	Any	<a href="https://www.ncbi.nlm.nih.gov/geo/">https://www.ncbi.nlm.nih.gov/geo/</a>
National Cancer Database	Any	<a href="https://www.facs.org/quality-programs/cancer/ncdb">https://www.facs.org/quality-programs/cancer/ncdb</a>
A public H&E-stained image dataset of colorectal cancer	Colorectal	<a href="https://zenodo.org/record/1214456">https://zenodo.org/record/1214456</a>
National Cancer Institute SEER Data	Any	<a href="https://seer.cancer.gov/data/">https://seer.cancer.gov/data/</a>
Genomics Data Commons	Any	<a href="https://portal.gdc.cancer.gov/">https://portal.gdc.cancer.gov/</a>
CPTAC Clear Cell Renal Cell Carcinoma (CCRCC) Discovery Study	Kidney cancer	<a href="https://wiki.cancerimagingarchive.net/display/Public/CPTAC-CCRCC">https://wiki.cancerimagingarchive.net/display/Public/CPTAC-CCRCC</a>
The Cancer Imaging Archive	Any	<a href="https://www.cancerimagingarchive.net/collections/">https://www.cancerimagingarchive.net/collections/</a>



# Table of content

- Project Overview
- Methodology
- Results
  - Skin and Breast Cancer
  - Kidney and Colon Cancer
  - **Spinal Cord, Cranial Nerves and Ovary**
  - Bronchus, Lung and Prostate Gland
- Next Steps

# Only about 37 publications related ML applications in Spinal Cord, Cranial Nerves and Ovary research were found due to rarity

## Research Overview

- Most of the data are coming from privately gathered information from different hospitals or medical center
- Despite the rarity of these cancer types, literature reviews show the potentials of machine learning techniques applied for both classification and regression prediction
- Machine learning maintains numerous advantages over conventional regression techniques, such as a reduced requirement for a priori knowledge on predictors and better ability to manage large datasets

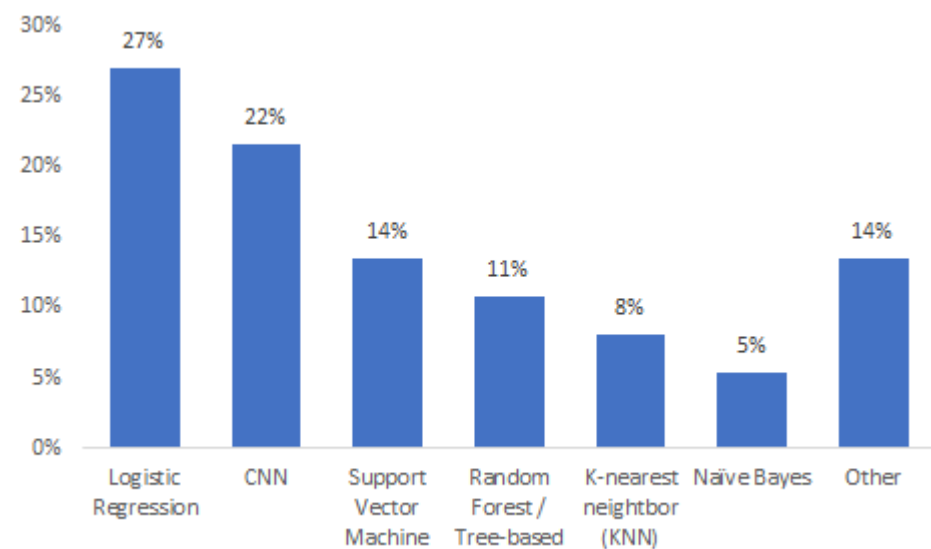
## Spinal Cord, Cranial Nerves and Ovary

### Use cases and objectives for ML research:

- Predictive Modeling (mostly classification):
  - Supplemented to clinical prognostic models for predicting the outcome of spinal cord injury (SCI)
  - Helped explore the non-linear relationships of patient features and prognostic functions
  - Classification in segmenting spinal cord and detecting ovarian cancer biomarker
  - Overperformed the traditional Logistic regression
- Pros for using ML:
  - Help describe nonlinear relationship
  - Less prior knowledge required
  - Improving diagnosis accuracy and efficiency
- Sample: 37

# The most common algorithms to use in Spinal Cord, Cranial Nerves and Ovary research are Logistic Regression and CNN, despite the rarity of datasets

**Algorithms Used for ML in Spinal Cord, Cranial Nerves and Ovary**



**Available Dataset**

Dataset	Cancer Type	Link
The Rick Hansen Spinal Cord Injury Registry (RHSCIR)	Spinal Cord	<a href="https://www.nature.com/articles/sc2011109">https://www.nature.com/articles/sc2011109</a>
Spinal Cord Injury Registry - NACTN	Spinal Cord	<a href="https://clinicaltrials.gov/ct2/show/results/NCT00178724">https://clinicaltrials.gov/ct2/show/results/NCT00178724</a>
European Multicentre Study of Human Spinal Cord Injury	Spinal Cord	<a href="https://clinicaltrials.gov/ct2/show/study/NCT01571531">https://clinicaltrials.gov/ct2/show/study/NCT01571531</a>
Ovary data	Ovary	<a href="https://figshare.com/articles/Raw_data/6025748">https://figshare.com/articles/Raw_data/6025748</a>

# Problems appeared in the research of Spinal Cord, Cranial Nerves and Ovary

- a. No public data source - Most of cases are privately collected from hospitals
- b. Ambiguous in describing the data samples / features
- c. Not robust performance due to the imbalance class distribution (less positive or diagnosed patients sample)
- a. Instability of algorithm feasibility when applying to different datasets
- a. Lack of automated tools for implementing or scaling the algorithm in real life

# Table of content

- Project Overview
- Methodology
- Results
  - Skin and Breast Cancer
  - Kidney and Colon Cancer
  - Spinal Cord, Cranial Nerves and Ovary
  - **Bronchus, Lung and Prostate Gland**
- Next Steps



# About 60 publications related machine learning applications in Lung, Bronchus and Prostate Cancer research were summarized

## Use cases and objectives for ML in cancer research

### Lung Cancer

- Predictive Modeling (mostly classification):
  - Patient survival probability
  - High risk-low risk classification of patients
  - Lung Cancer treatment response
- Deep Learning:
  - Image classification for lung cancer prognosis
  - Effectiveness of therapeutic antibody targeting in the body

### Bronchus Cancer

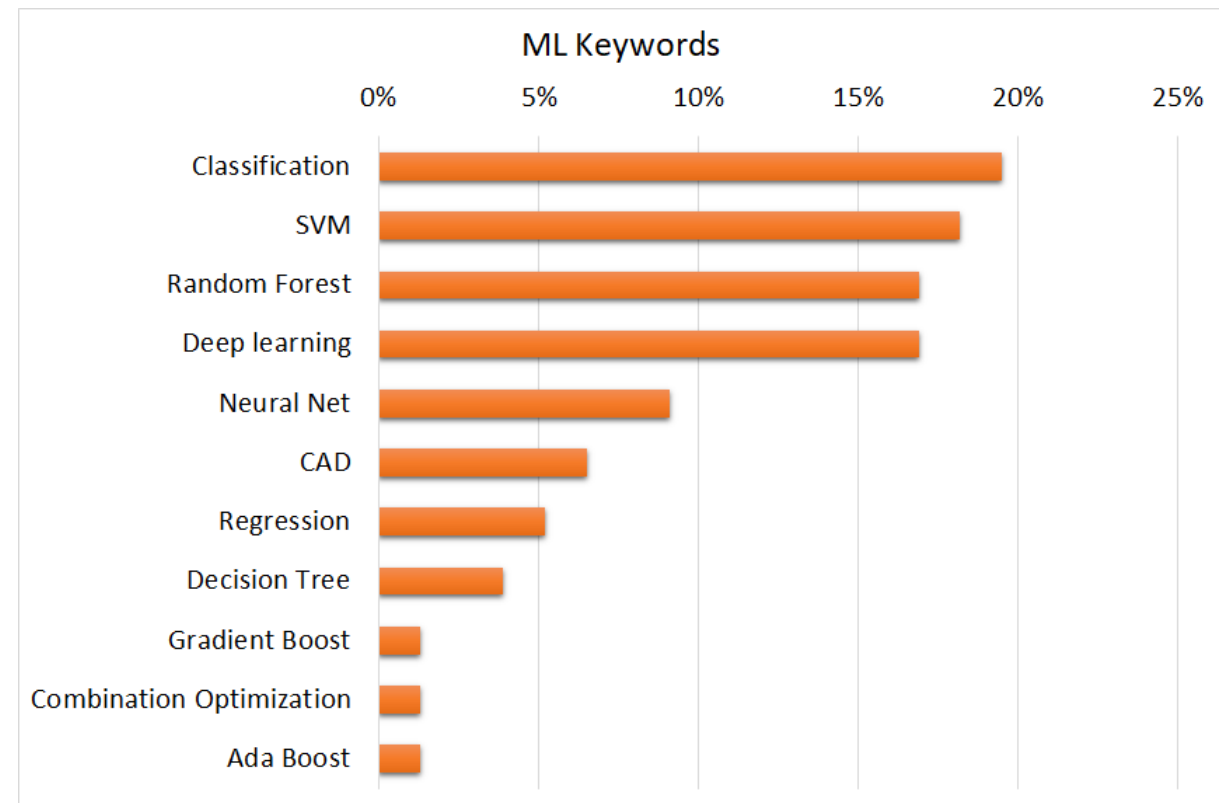
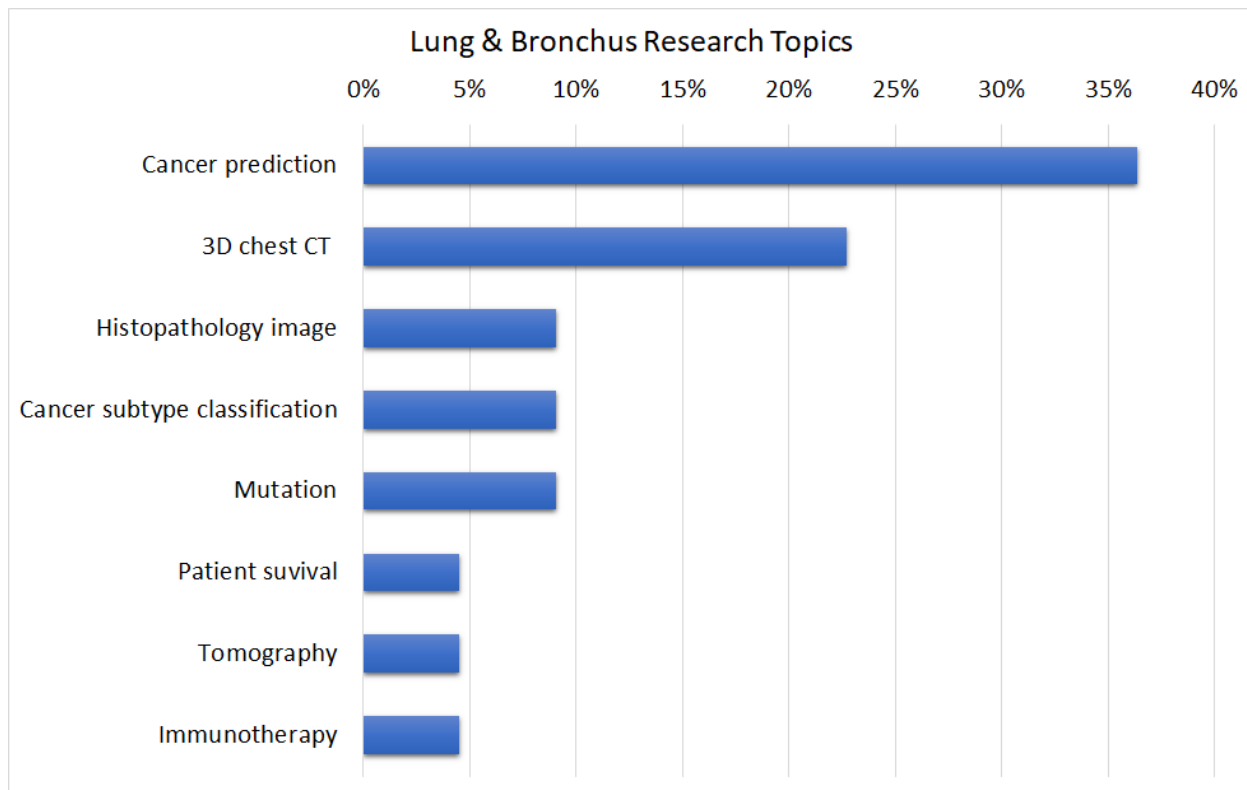
- Predictive Modeling (mostly classification):
  - Predict treatment outcomes of chemotherapy
- Deep Learning:
  - Automated anatomical labeling of bronchial branches extracted from 3D CT images
- Limited research as compared to lung and prostate cancer

### Prostate Cancer

- Predictive Modeling (mostly classification):
  - Prostate cancer probability prediction for early diagnosis
  - Survival probability
- Deep Learning:
  - Automated gleason grading of prostate cancer tissue
  - Prostate cancer magnetic resonance imaging

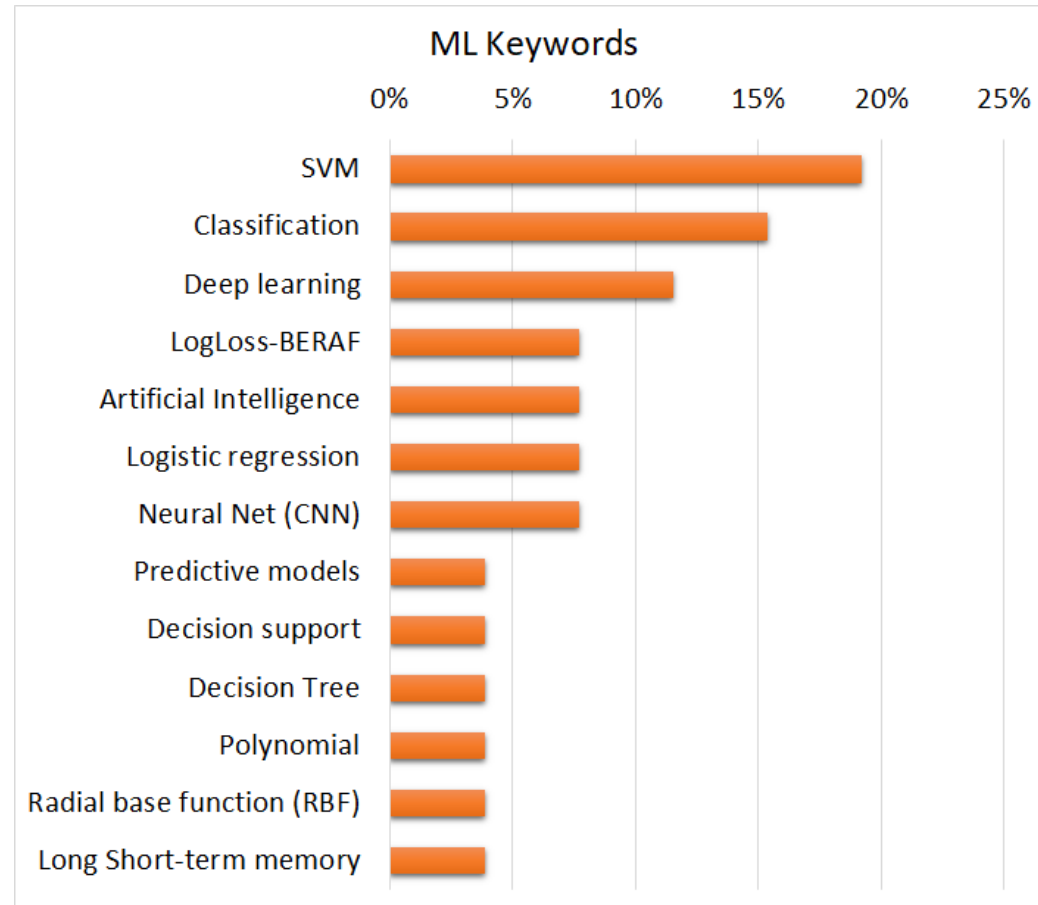
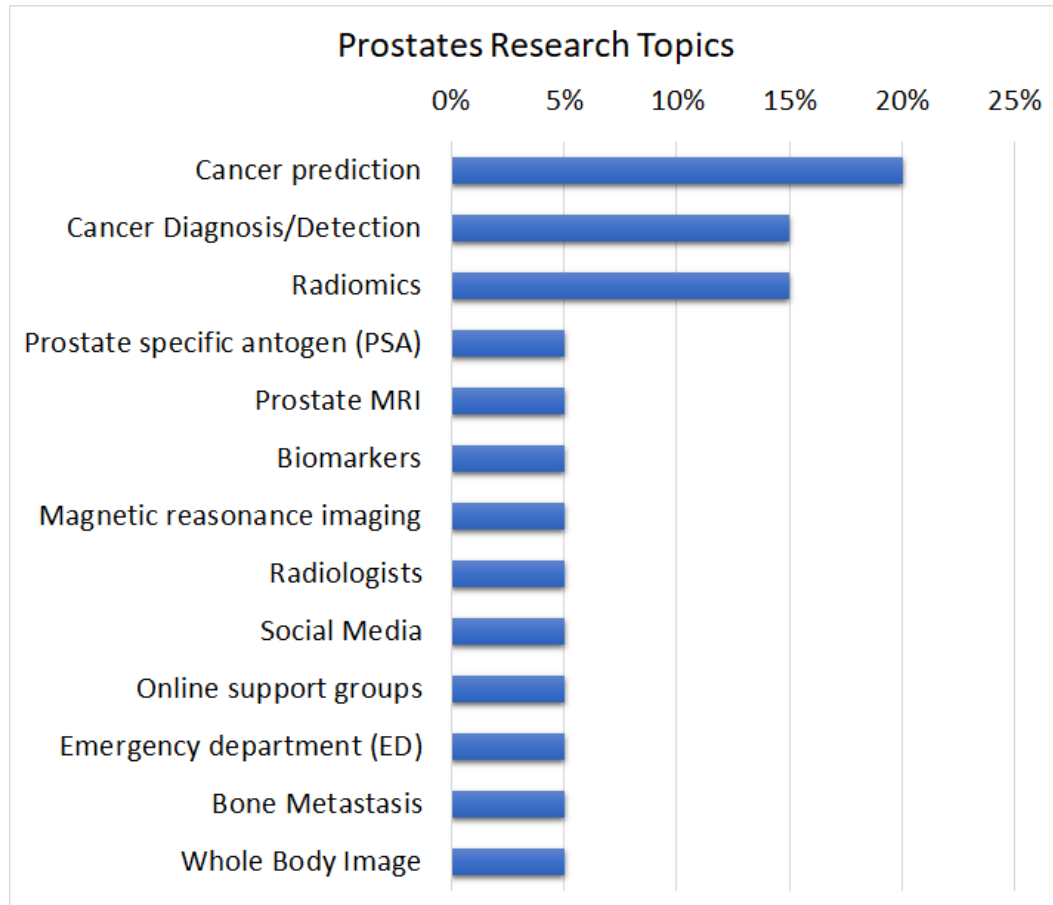
# Cancer Prediction is the most popular research topic and Classification is the most frequently used machine learning technique

## Top research topics and machine learning key words..



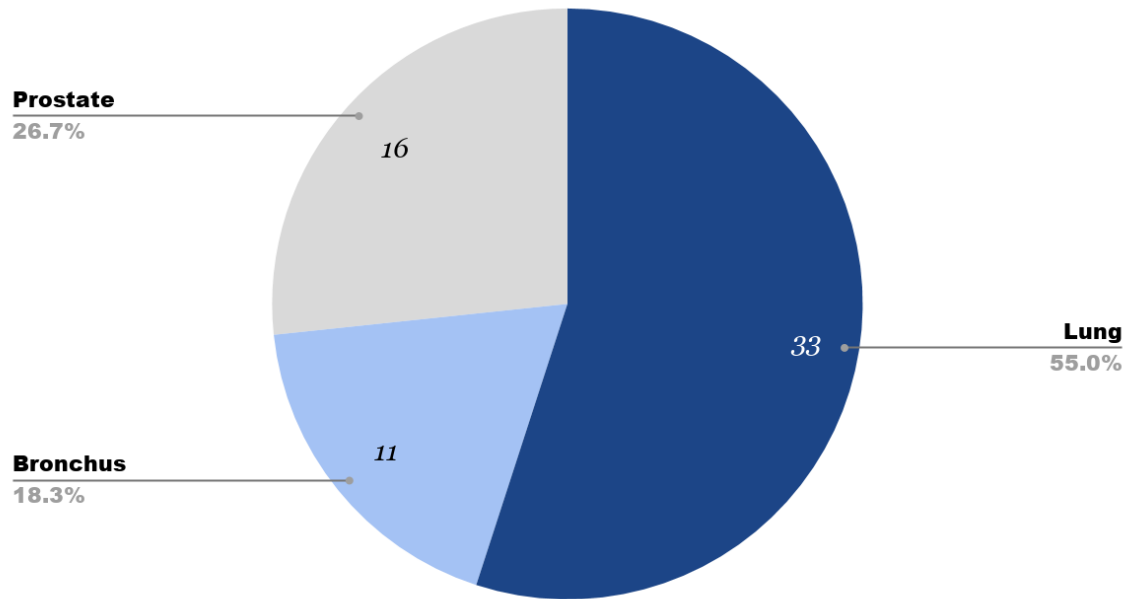
# Cancer Prediction is the most popular research topic and SVM is the most frequently used machine learning technique

## Top research topics and machine learning key words..

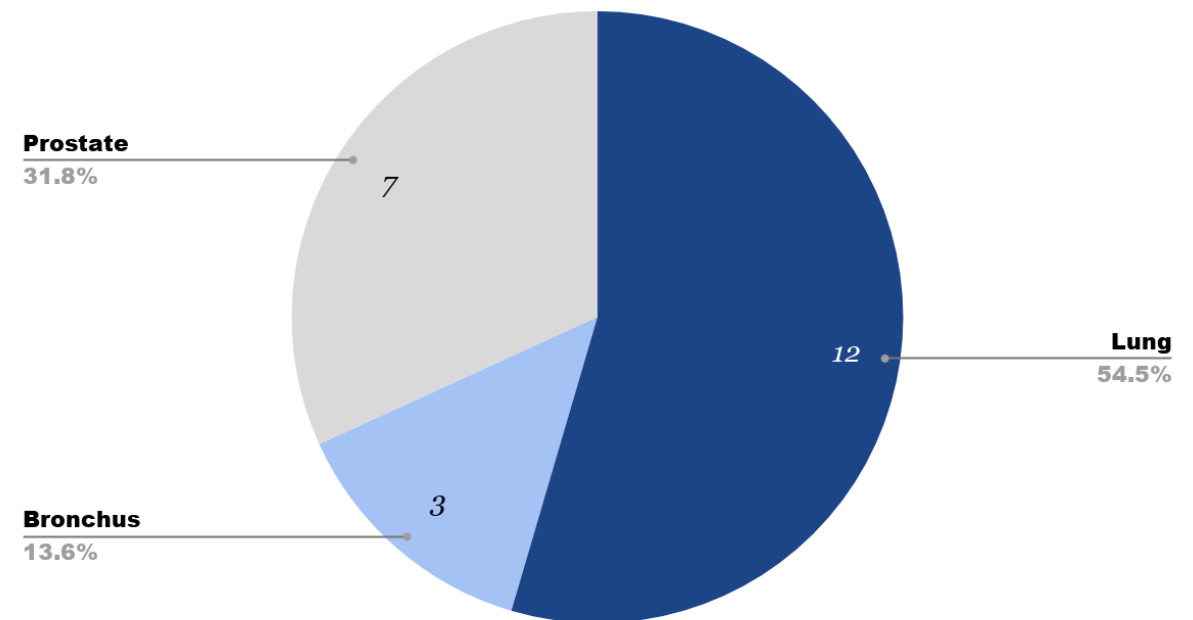


# Gathered data sources for ~37% publications summarized and most of them are publicly available

Number of publications summarized by type of cancer



Number of datasets identified by type of cancer



In total 22 datasets were collected in total out of 60 papers summarized for Prostate, Bronchus and Lung cancer

# 20 Datasets were identified from Lung and Prostate cancer research publications

S.No.	Dataset (L - Lung, P - Prostate)	Link	S.No.	Dataset (L - Lung, P - Prostate)	Link
1	SEER Incidence Data - L	<a href="https://seer.cancer.gov/data/">https://seer.cancer.gov/data/</a>	11	Biomarker Data - L	<a href="https://www.iprox.org/page/project.html?id=IPX0001153000">https://www.iprox.org/page/project.html?id=IPX0001153000</a>
2	Genomic Data Commons - L	<a href="https://portal.gdc.cancer.gov/">https://portal.gdc.cancer.gov/</a>	12	Immunotherapy Data - L	<a href="https://wiki.cancerimagingarchive.net/display/Public/">https://wiki.cancerimagingarchive.net/display/Public/</a>
3	LIDC-IDRI - L	<a href="https://wiki.cancerimagingarchive.net/display/Public/LIDC-IDRI">https://wiki.cancerimagingarchive.net/display/Public/LIDC-IDRI</a>	13	Biomarker Data - P	<a href="https://luisrueda.myweb.cs.uwindsor.ca/datasets/Hierarchical-Prostate-Cancer-Gleason.rar">https://luisrueda.myweb.cs.uwindsor.ca/datasets/Hierarchical-Prostate-Cancer-Gleason.rar</a>
4	NLST Pathology Images - L	<a href="https://cdas.cancer.gov/learn/nlst/images/">https://cdas.cancer.gov/learn/nlst/images/</a>	14	OSG Data - P	<a href="https://github.com/tharindurb/PRIME">https://github.com/tharindurb/PRIME</a>
5	NSCLC-Radiomics - L	<a href="https://wiki.cancerimagingarchive.net/display/Public/NSCLC-Radiomics">https://wiki.cancerimagingarchive.net/display/Public/NSCLC-Radiomics</a>	15	GEO Dataset - P	<a href="https://www.ncbi.nlm.nih.gov/gds">https://www.ncbi.nlm.nih.gov/gds</a>
6	NSCLC-Radiomics-Genomics - L	<a href="https://wiki.cancerimagingarchive.net/display/Public/NSCLC-Radiomics-Genomics">https://wiki.cancerimagingarchive.net/display/Public/NSCLC-Radiomics-Genomics</a>	16	Cancer Genome Atlas - P	<a href="https://docs.google.com/spreadsheets/d/12CXisExMibP_11GJNZp0Uit7-1VX92_eK0GfWysbmPK/edit#gid=201540152">https://docs.google.com/spreadsheets/d/12CXisExMibP_11GJNZp0Uit7-1VX92_eK0GfWysbmPK/edit#gid=201540152</a>
7	Gene-expression Data - L	<a href="https://www.ncbi.nlm.nih.gov/geo/query/acc.cgi?acc=GSE58661">https://www.ncbi.nlm.nih.gov/geo/query/acc.cgi?acc=GSE58661</a>	17	Radiographs - P	<a href="https://figshare.com/articles/Data_from_An_Invstigatoin_of_Machine_Learning_Methods_in_Delta-radiomics_Feature_Analysis/9943334">https://figshare.com/articles/Data_from_An_Invstigatoin_of_Machine_Learning_Methods_in_Delta-radiomics_Feature_Analysis/9943334</a>
8	RIDER Collections - L	<a href="https://wiki.cancerimagingarchive.net/display/Public/RIDER+Collections">https://wiki.cancerimagingarchive.net/display/Public/RIDER+Collections</a>	18	UTI Diagnosis Data - P	<a href="https://doi.org/10.1371/journal.pone.0194085.s001">https://doi.org/10.1371/journal.pone.0194085.s001</a>
9	TCGA Dataset - L	<a href="https://gdc.cancer.gov/about-data/publications/pancanatlas">https://gdc.cancer.gov/about-data/publications/pancanatlas</a>	19	CPC-GENE Data - P	<a href="http://www.cbioportal.org/study/summary?id=prad_cpcg_2017">http://www.cbioportal.org/study/summary?id=prad_cpcg_2017</a>
10	NGC Data - L	<a href="https://stm.sciencemag.org/content/10/457/eaar7939/tab-figures-data">https://stm.sciencemag.org/content/10/457/eaar7939/tab-figures-data</a>	20	PRAD-FR - P	<a href="https://dcc.icgc.org/projects/PRAD-FR">https://dcc.icgc.org/projects/PRAD-FR</a>



# Table of content

- Project Overview
- Methodology
- Results
  - Skin and Breast Cancer
  - Kidney and Colon Cancer
  - Spinal Cord, Cranial Nerves and Ovary
  - Bronchus, Lung and Prostate Gland
- **Summary and Next Steps**

# Image processing and predictive modeling are common machine learning applications in cancer research

There are some commonalities in the research carried out across different types of cancer

## **Image processing for cancer detection**

- Machine learning based image processing techniques to identify the cancer type and further classify patients as high risk and low risk groups.
- Deep Learning and Convolutional Neural Networks (CNNs) are most common techniques followed by SVM

## **ML applications in cancer prognosis**

- Predictive modeling to improve accuracy in predicting cancer susceptibility, recurrence and survival prediction.
- However, there is lack of external validation or testing regarding the performance of these predictive models

# Future Work

- The datasets gathered can be further used to conduct research at FNL
- Building on the datasets identified, one can dive into those datasets to validate the integrity of the research publications
- Continue to gather more datasets available, as are a still lot research papers available for some cancer areas that are not summarized during this project
- In future, loop in more senior people in emails to authors to improve response rate
- NLP techniques can be applied on research summary table (e.g. to identify popular softwares used, ML packages etc.)
- Conduct a trend analysis, which ML techniques are gaining popularity in the recent years



# 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 medical science. We've learned a great deal and hope to be able to work with FNL again.

We have identified *79 publicly available datasets* for machine learning publications in the oncology space.

*Any questions?*