



The Elizabeth H.  
and James S. McDonnell III

**McDONNELL  
GENOME INSTITUTE**  
at Washington University



@obigriffith

griffithlab.org

# Artificial Intelligence in Genomic Medicine

Obi L. Griffith <obigriffith@wustl.edu>

# What is Artificial Intelligence?

---

- Artificial Intelligence (AI) is the theory and development of computer systems able to perform tasks that normally require human intelligence, such as visual perception, speech recognition, decision-making, and other complex tasks



# Examples of advanced AI appearing throughout society

- Manufacturing robots
- Automated financial investing
- Social media monitoring
- Smart assistants
- Conversational bots
  - ChatGPT



Many others...



# Applications in biomedical field far too many to list

- Medical image analysis
  - MRI segmentation
  - Histopathological cancer classification
  - Mammographic lesion detection
- Protein structure prediction
- Genome analysis
  - Gene expression inference
  - Enhancer prediction
  - Variant pathogenicity prediction

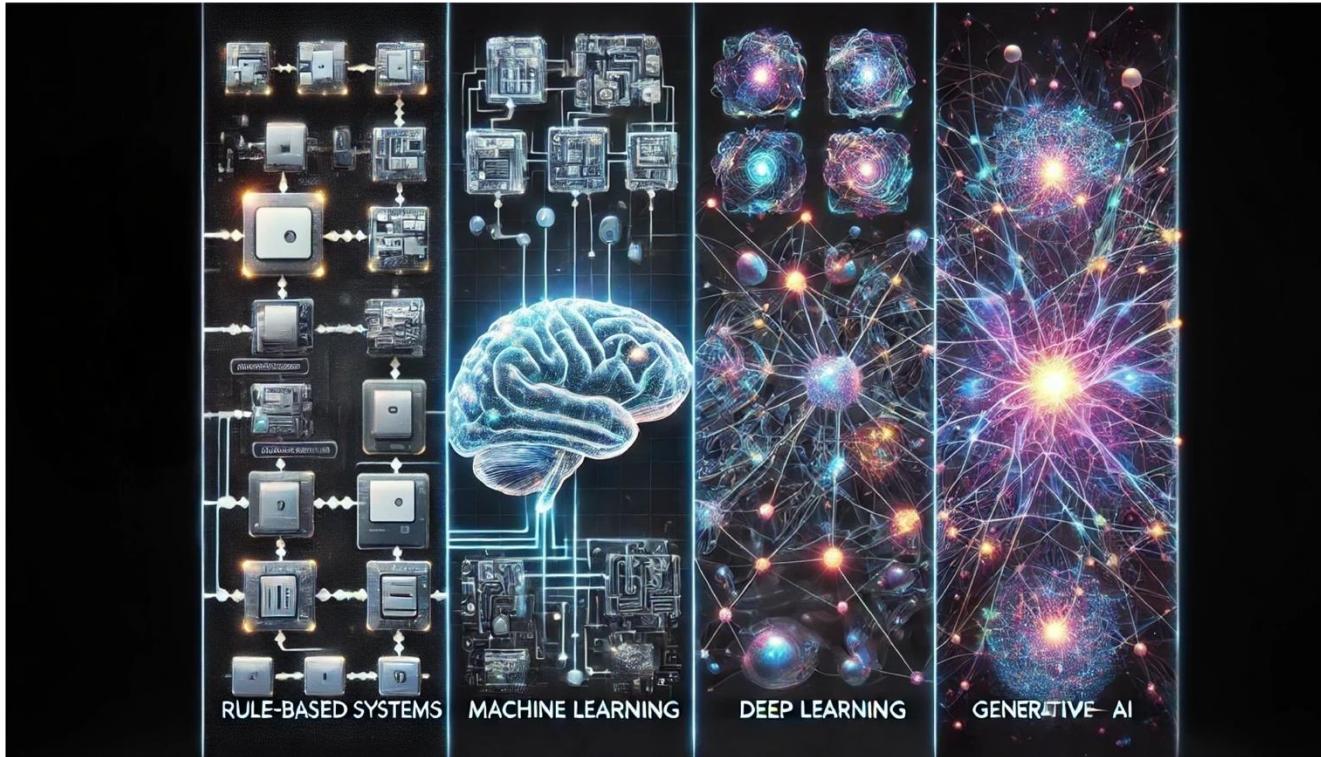


# What are the key challenges of AI?

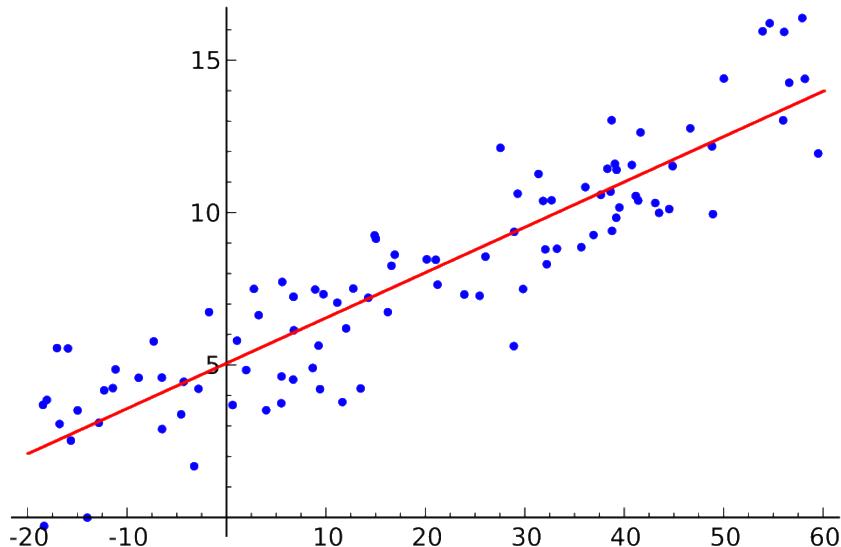
- Key Challenges of AI: knowledge representation, natural language processing (NLP), learning, perception, planning, reasoning, motion, social intelligence and general intelligence
  - Knowledge representation - catalog concepts and relationships using knowledgebases, ontologies, and structured data, for computational use
  - Natural language processing (e.g., text-mining) - allow computers to read and understand human language
  - Learning (e.g., machine learning) - develop computer algorithms that improve automatically through experience
    - Unsupervised - Find novel patterns in data without human input or labels
    - Supervised - Predict unknowns from labeled knowns (e.g., classifiers)



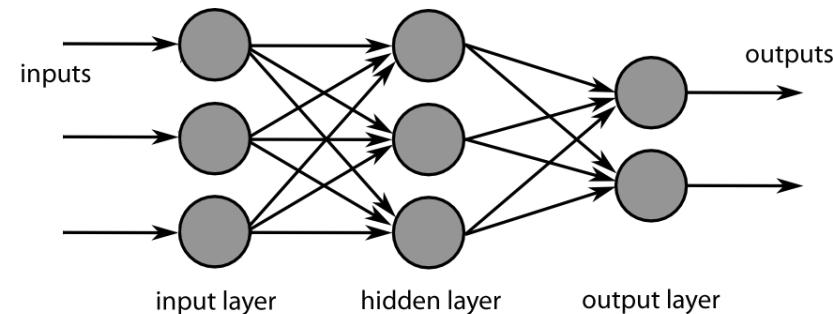
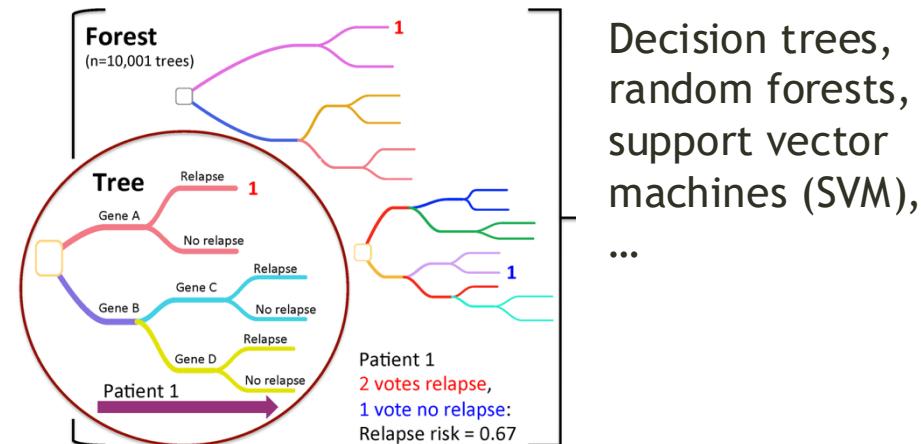
# AI has gone through several eras



# Learning -> Machine Learning -> Deep Learning



Linear regression



multi-layered artificial neural networks



# General considerations

- Establish training and testing datasets
  - Are there batch effects to consider?
  - Randomly split data
    - Make sure balanced for key features (sex, target class, etc)
  - Identify independent test sets
- Define target classes (e.g., Relapse vs No Relapse)
  - How is relapse defined?
  - Are the classes equal sizes? If not, how to mitigate?
- Are there features which should be filtered out?
- Are there new features we wish to design/engineer?
- Is the model overfit/overtrained?



# Let's walk through an example

Griffith et al. *Genome Medicine* 2013, **5**:92  
<http://genomemedicine.com/content/5/10/92>



RESEARCH

Open Access

## A robust prognostic signature for hormone-positive node-negative breast cancer

Obi L Griffith<sup>1,2\*</sup>, François Pepin<sup>1,3</sup>, Oana M Enache<sup>1</sup>, Laura M Heiser<sup>1,4</sup>, Eric A Collisson<sup>5</sup>, Paul T Spellman<sup>1,6\*</sup>  
and Joe W Gray<sup>1,4\*</sup>



<https://www.biostars.org/p/85124/>

Nick Spies



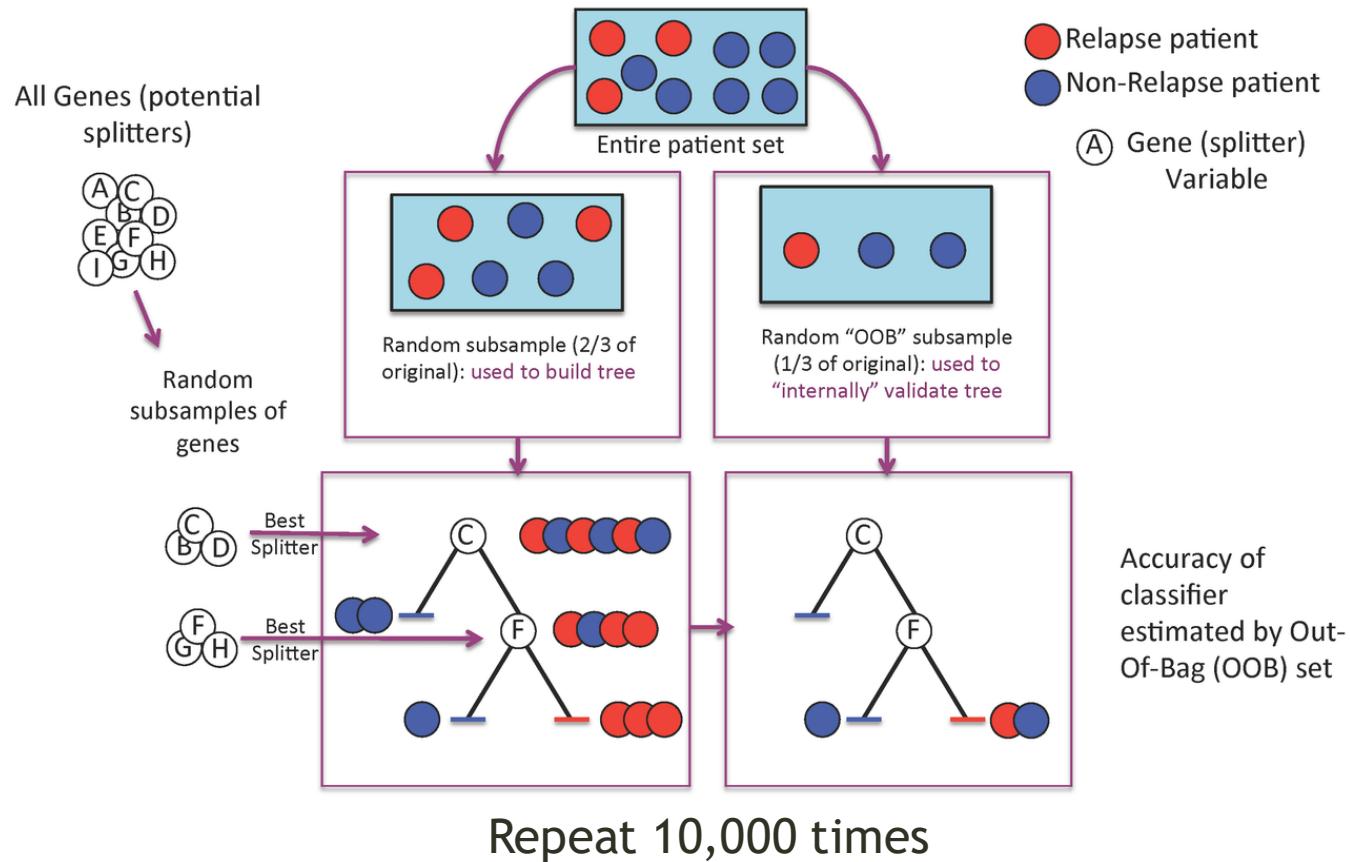
# Let's walk through an example

- Let's build a classifier of relapse in ER+ breast cancer
- Suppose we have a gene expression dataset for a set of patients, with known relapse status

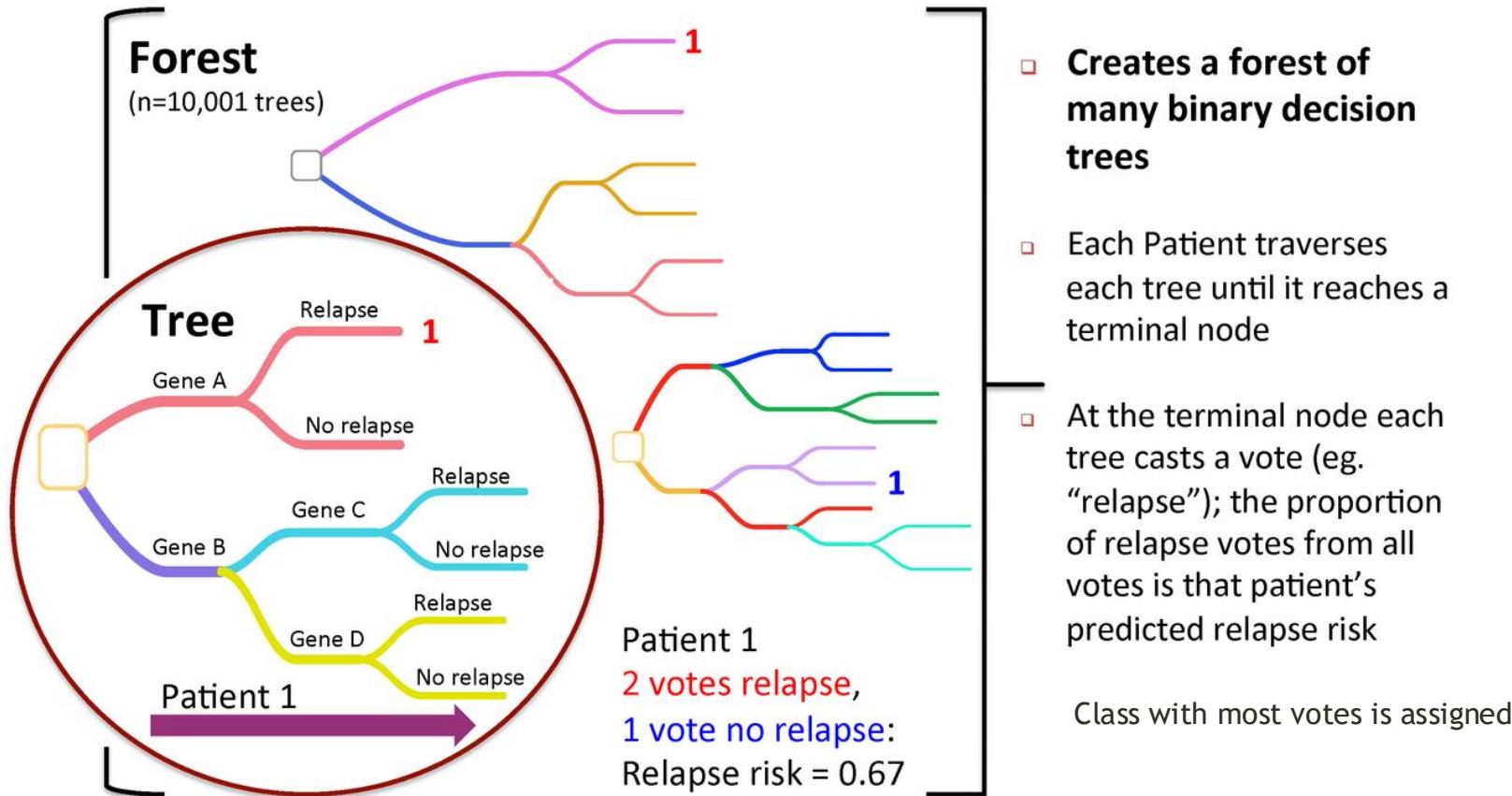
	Relapse	Relapse	No Relapse	No Relapse
	Patient1	Patient2	Patient3	Patient4
GeneA	4.7	5.2	1.1	0.9
GeneB	0.5	0.7	6.8	7.1
GeneC	2.3	3.1	2.4	3.2



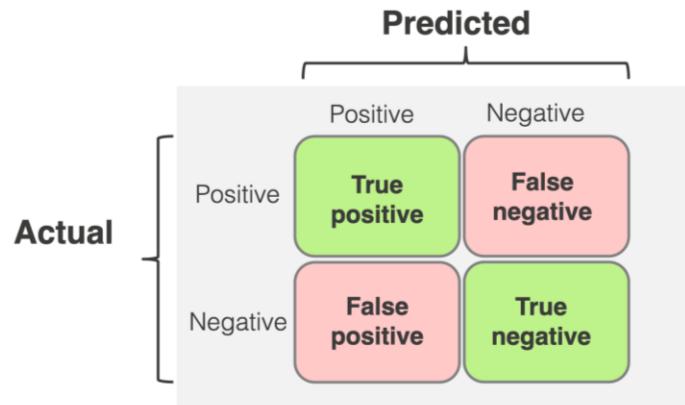
# How do you build (train) a Random Forest model?



# How do you predict with a Random Forest model?



# How does our model perform? Let's look at the dreaded confusion matrix



Actual	Predicted	
	Relapse	NoRelapse
Relapse	49	58
NoRelapse	39	140

**Sensitivity/Recall** (true positive rate) is how well a test can identify true positives

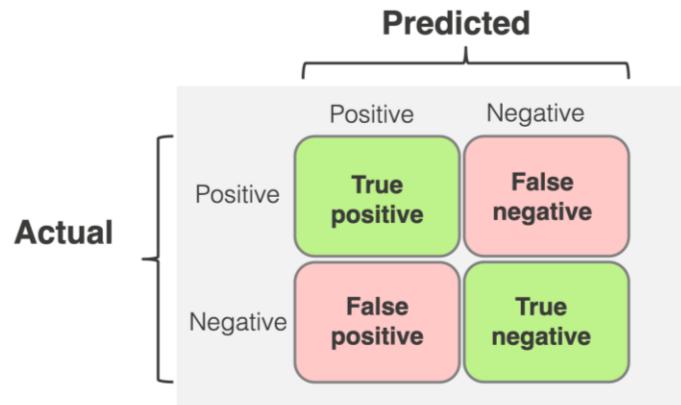
$$49 / (49+58) = 0.458$$

**Specificity>Selectivity** (true negative rate) is how well a test can identify true negatives

$$140 / (39+140) = 0.782$$



# How does our model perform? Let's look at the dreaded confusion matrix



Actual	Predicted	
	Relapse	NoRelapse
Relapse	49	58
NoRelapse	39	140

**Relapse** class error rate is how many true Relapses were miscalled

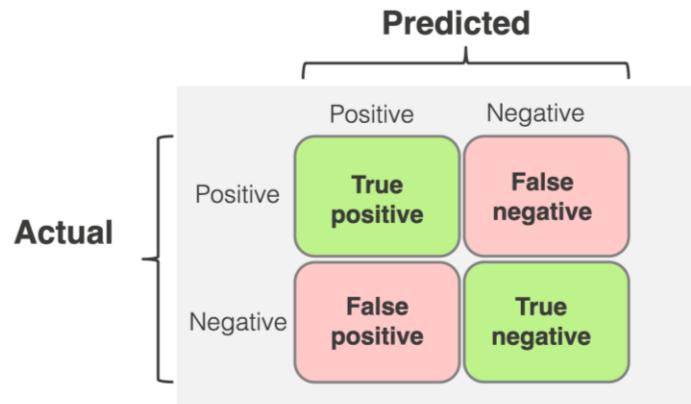
$$58 / (49+58) = 0.542$$

**No Relapse** class error rate is how many true No Relapses were miscalled

$$39 / (39+140) = 0.218$$



# How does our model perform? Let's look at the dreaded confusion matrix



Actual	Predicted	
	Relapse	NoRelapse
Relapse	49	58
NoRelapse	39	140

Overall accuracy

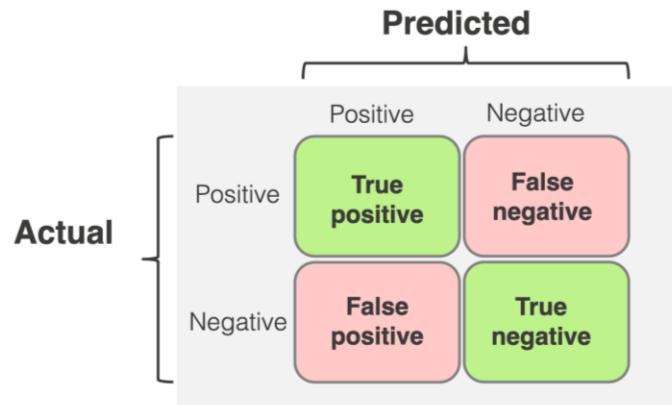
$$(49+140)/(49+58+39+140) = 0.660$$

Overall error

$$(39+58)/(49+58+39+140) = 0.340$$



# How does our model perform? Let's look at the dreaded confusion matrix



Actual	Predicted	
	Relapse	NoRelapse
Relapse	49	58
NoRelapse	39	140

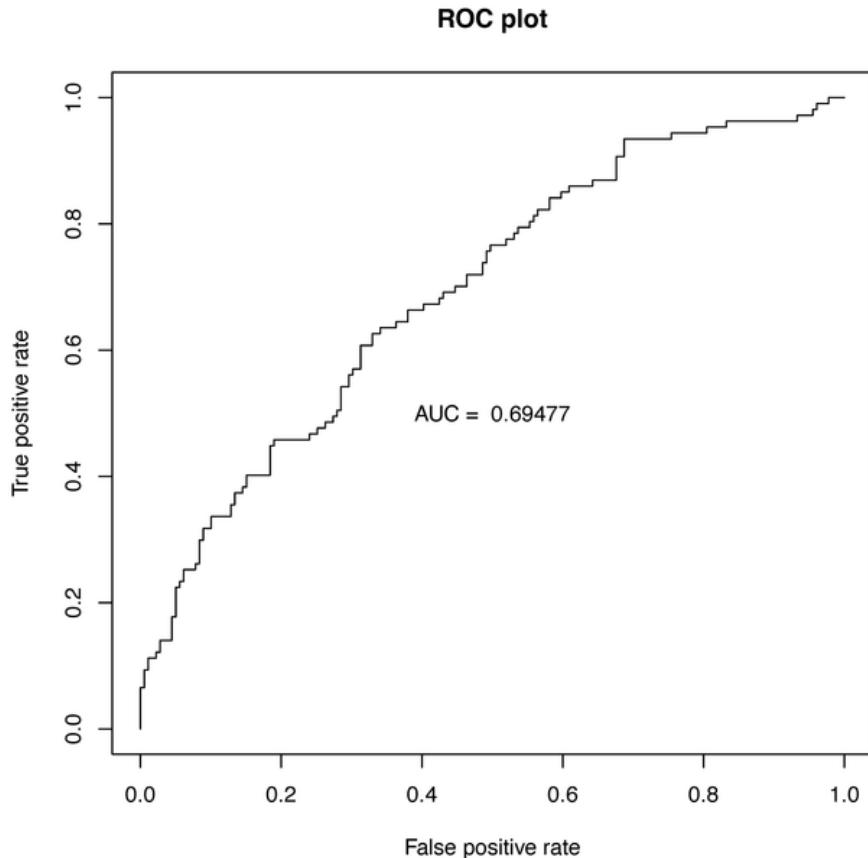
Precision is the share of true positive predictions in all positive predictions

$$49 / (49 + 39) = 0.557$$

Other popular metrics: F1 score, PPV, NPV, etc



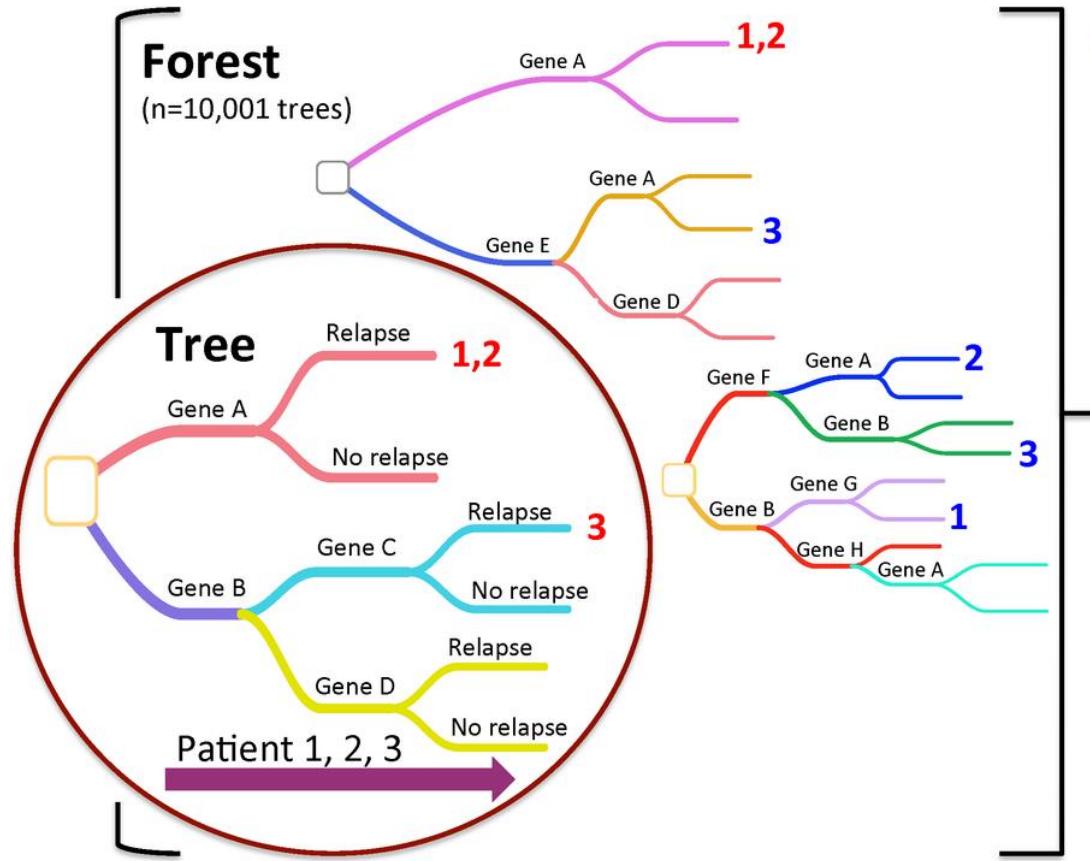
# Receiver Operating Characteristic (ROC) plots TPR (Sensitivity) vs FPR for all possible model thresholds



The ROC AUC summarizes how well the model can distinguish between positive and negative classes



# How do we know which features are important in our model

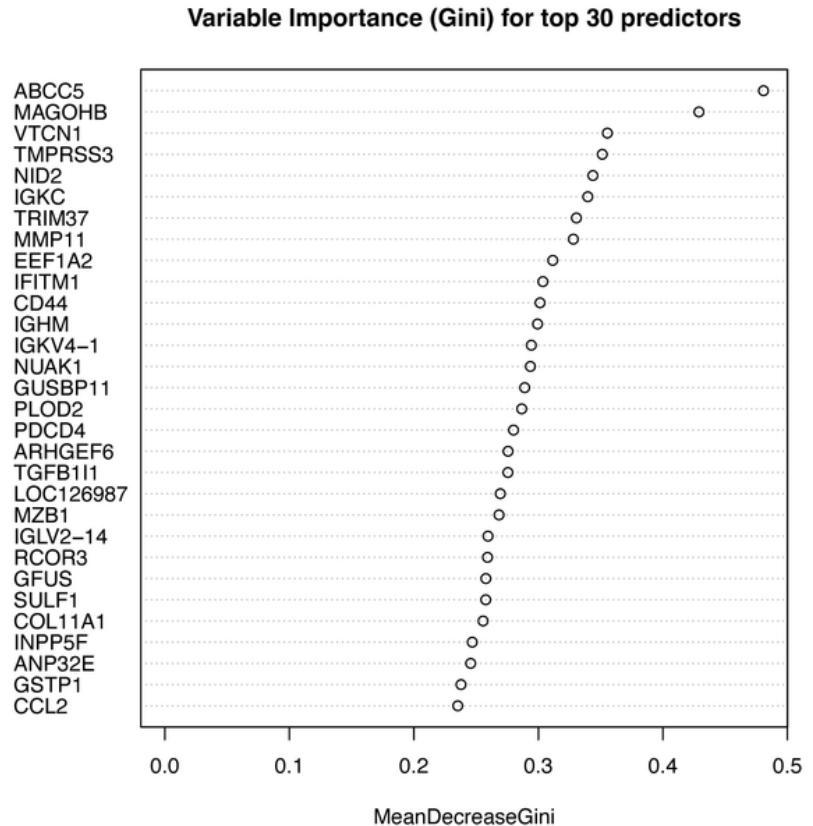


- The more often a gene is chosen as a splitter variable, the higher its “Variable Importance” – This can be used to prioritize which genes to select for an assay with limited gene measurements

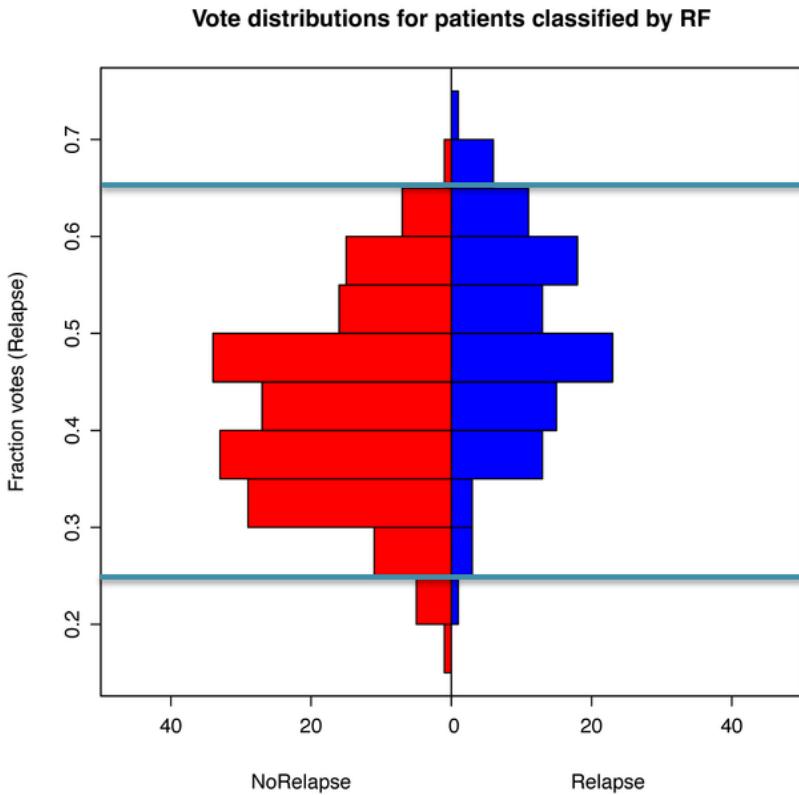
Gene	Var. Imp.
Gene A	0.67
Gene B	0.20
Gene D	0.13
...	...



# Variable importance could be used to prioritize genes for a targeted assay



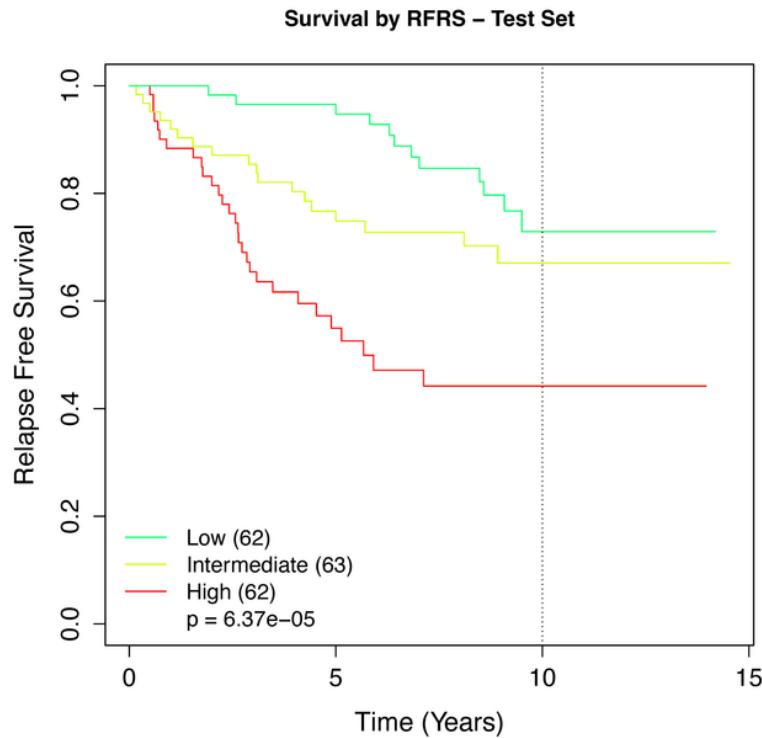
# Even a model with poor overall performance may be very accurate for a subset of cases



high-risk,  
intermediate-risk and  
low-risk groups can be  
defined



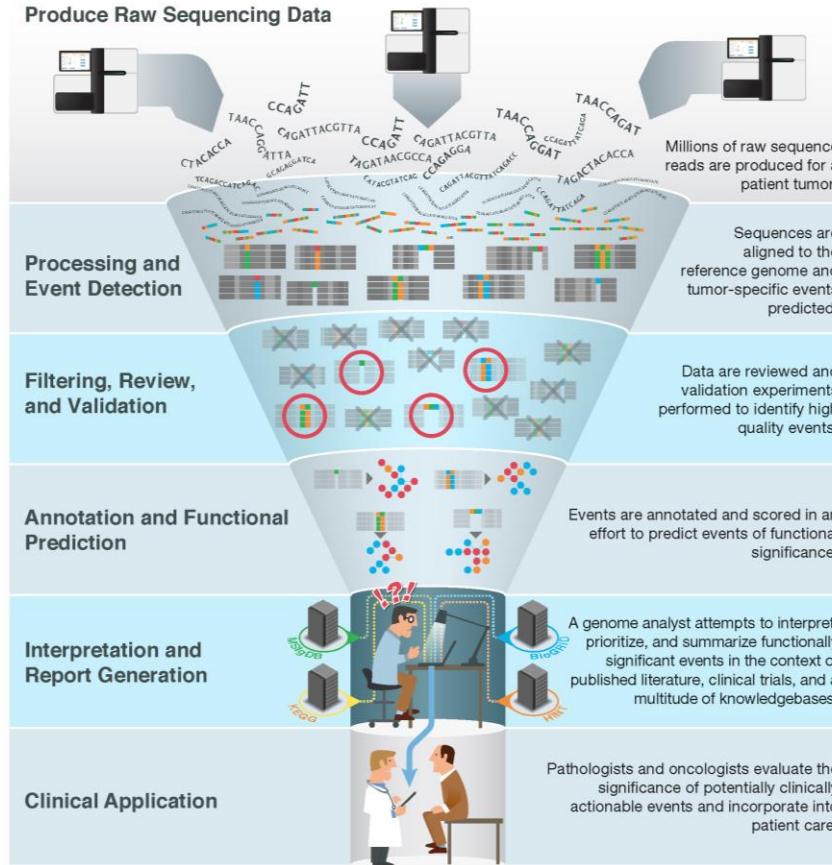
# These risk groups can be related back to survival



This is how essentially how Oncotype DX, MammaPrint, and ProSigna assays were developed - used routinely in breast cancer care



# How else can AI be applied to genomic medicine?



Variant calling

Manual review

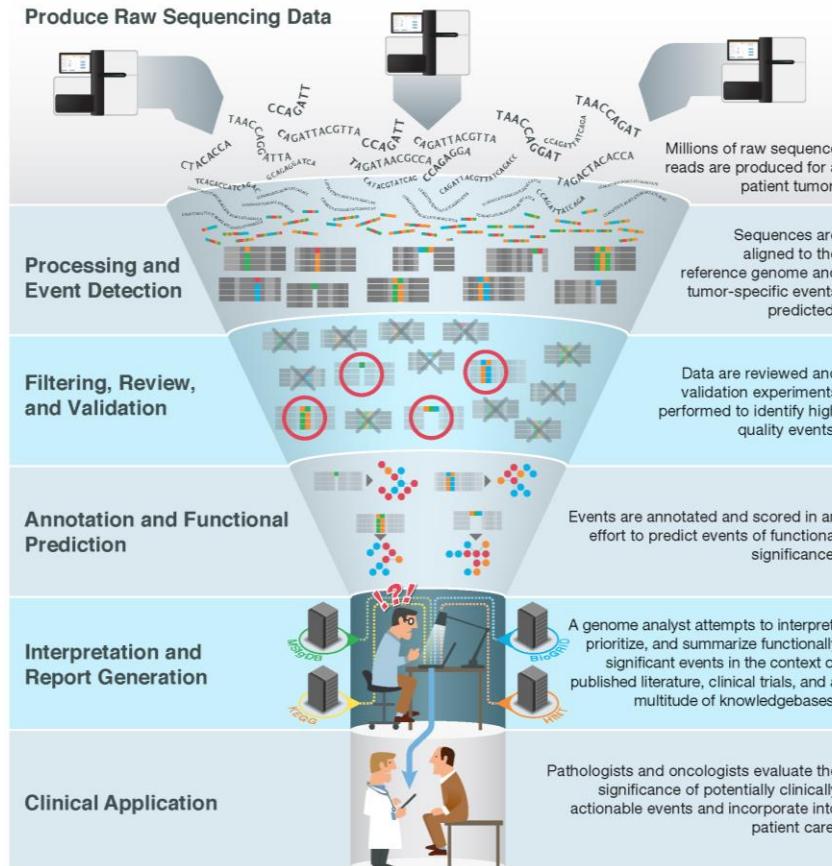
Functional prediction

Clinical Relevance

Decision-making



# How else can AI be applied to genomic medicine?



Variant calling

Manual review

Functional prediction

Clinical Relevance

Decision-making



# I) Manual review - how do we know a cancer variant is real?

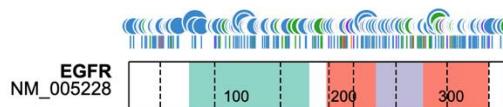


# Variant calling from raw sequence data is a critical step in most cancer genomics research and clinical laboratory genomics workflows

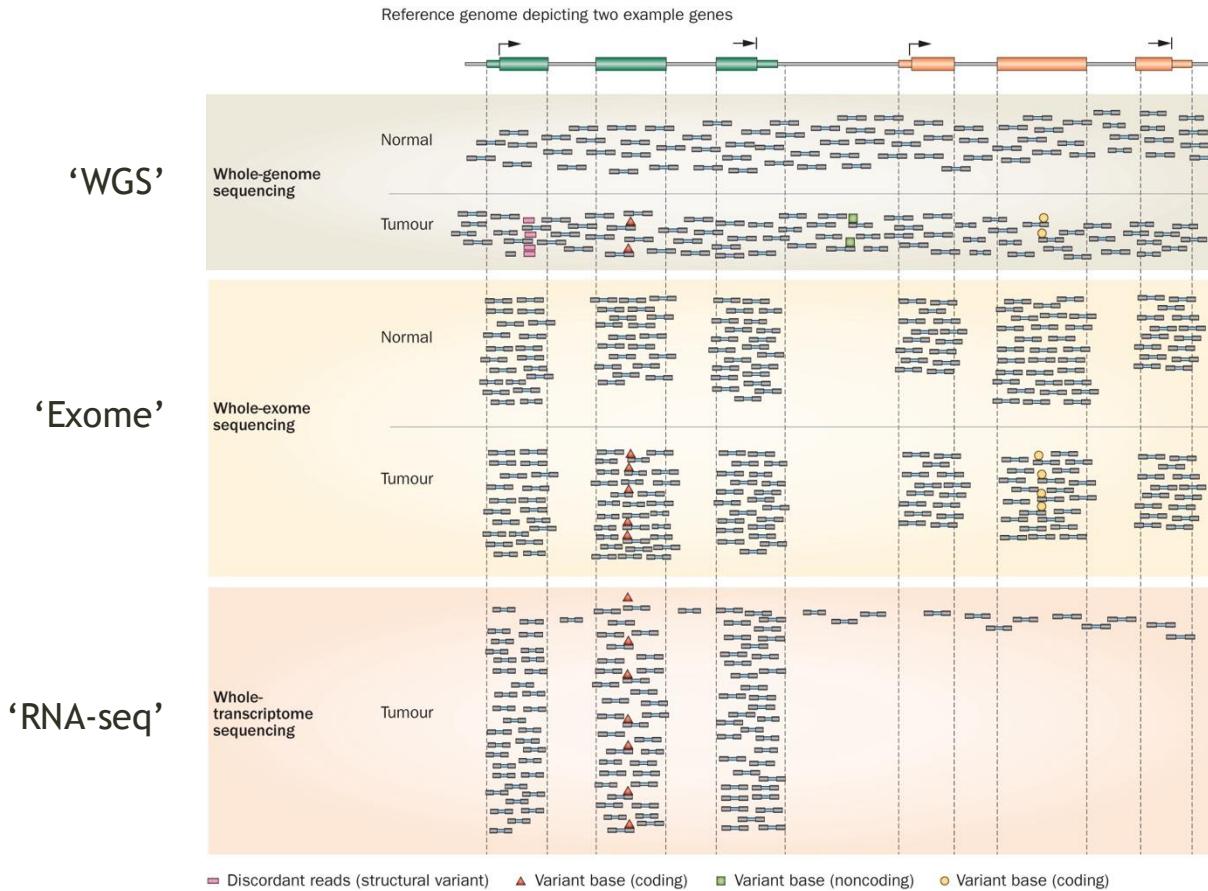


ABOUT

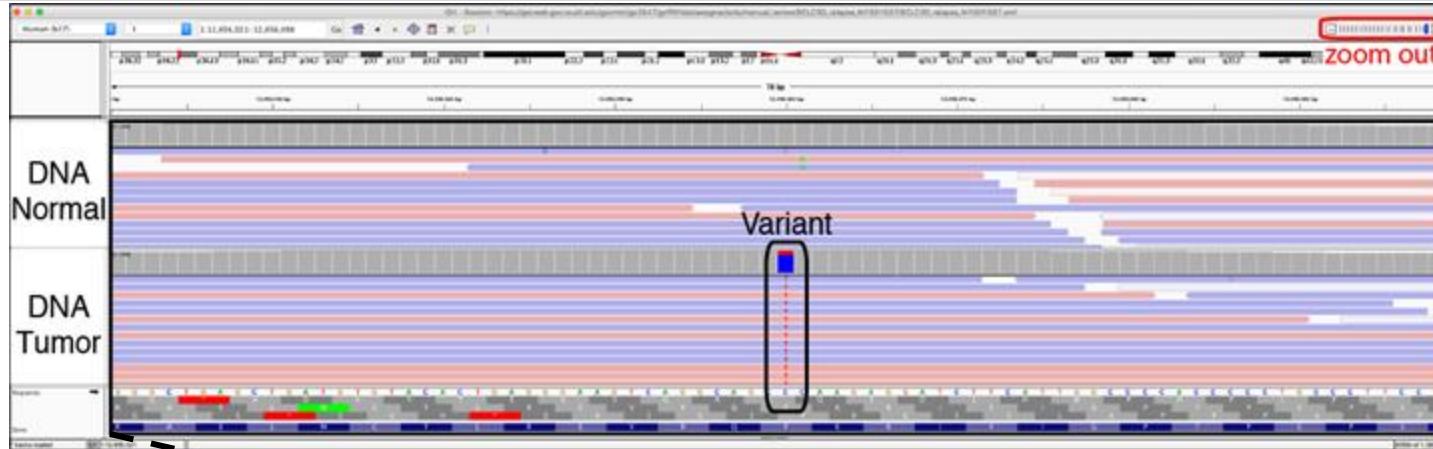
GENOMIC ALTERATIONS	
GENE	INTERPRETATION
ALTERATION	
PIK3CA H1047R	Mutations in PIK3CA have been reported in 26% to 33% of breast cancer cases (COSMIC, Jun 2012 and Kalinsky et al., 2009; 19671852). Activating mutations in PIK3CA, such as the one seen here, may predict sensitivity to inhibitors of PI3 kinase or its downstream signaling pathway (the PI3K/Akt/mTOR pathway) (Huang et al., 2007; 18079394). The mTOR inhibitors temsirolimus and everolimus have been tested in several clinical trials in breast cancer, and have been approved by the FDA for use in other tumor types. Inhibitors of PI3K and Akt are currently in clinical trials in breast cancer, alone or in combination with other therapies. PIK3CA mutations may play a role in resistance to hormonal therapy in ER+ breast cancers (Miller et al., 2011; 22114931). Activating mutations in PIK3CA may also confer resistance to anti-Her2 therapies (Chakrabarty et al., 2010; 20581867, Kataoka et al., 2010; 19633047, Wang et al., 2011; 21676217); combined inhibition of Her2 and the PI3K pathway may be required in tumors with ERBB2 amplification and PIK3CA mutation, though this remains an area of active investigation.
CCND1 amplification	CCND1 amplification has been reported in approximately 10-15% of invasive breast cancers, more frequently in BRCA-negative cancers (Elshiekh et al., 2008; 17653856, Bane et al., 2011; 21327470). There are no approved therapies that directly target the protein product of CCND1 (Cyclin D1); however, CCND1 amplification may predict sensitivity to inhibitors of Cdk4 and Cdk6, which are currently under investigation in clinical trials. Overexpression of Cyclin D1 has also been associated with resistance to endocrine therapy in breast cancer (reviewed in Lange et al., 2011; 21613412; Musgrove and Sutherland, 2009; 19701242, Butt et al., 2005; 16113099).
CDH1 E167*	CDH1 mutations are present in approximately 17% of breast cancers, and more often in luminal type cancers (COSMIC, Jun 2012, Hollestelle et al., 2010; 19593635). Loss of the E-cadherin protein, which is encoded by the CDH1 gene, has been associated with poor prognosis in triple negative breast cancer (Kashiwagi et al., 2010; 20551954, Tang et al., 2011; 21519872). Presently, there are no targeted therapies to address loss of CDH1/E-cadherin.



# Whole genome, exome and transcriptome sequencing allows us to detect and confirm many different types of ‘omic events



# Manual review involves human inspection of sequence alignments to identify artifacts missed by automated variant callers



Literature is full of references to manual review and ad hoc filtering

"Mutations...were called with MuTect and filtered with oxidation and panel of normal samples filters to remove artefacts..."



## Recurrent and functional regulatory mutations in breast cancer

Esther Rhee-Shay<sup>1,2</sup>, Premaan Parasuraman<sup>3</sup>, Sunita Grimaldy<sup>4</sup>, Grace Tsui<sup>5</sup>, Jessie M. Engqvist<sup>1,2</sup>, Jangki Kim<sup>1</sup>, Michael S. Lawrence<sup>1</sup>, Aman Taylor-Weiner<sup>1</sup>, Steven Rodriguez-Carmon<sup>1</sup>, Maya Rabinowitz<sup>1</sup>, Julian Henn<sup>1</sup>, Chip Stewart<sup>1</sup>, Vasil E. Marinovik<sup>1</sup>, Peter Barjaktov<sup>1</sup>, Maria L. Cortes<sup>1</sup>, Saia Sepehri<sup>1</sup>, Gerae Cebulski<sup>1</sup>, Adam Tracy<sup>1</sup>, Trevor J. Pugh<sup>1</sup>, Jessie Lee<sup>1</sup>, Zongqi Zheng<sup>1</sup>, Leaf W. Elliston<sup>1</sup>, A. John Sahraian<sup>1</sup>, Jessie S. Boddie<sup>1</sup>, Stacey B. Gabriel<sup>1</sup>, Matthew Maynor<sup>1</sup>, Todd R. Golich<sup>1</sup>, Jose Basulto<sup>1</sup>, Alfredo Hafalero Minaya<sup>1</sup>, Tricia Shobola<sup>1</sup>, Andrei Bernstein<sup>1</sup>, Eric L. Landner<sup>1</sup> and Ged Getz<sup>1,2,3,4,5</sup>



## An immunogenic personal neoantigen vaccine for patients with melanoma

Pannik, A., et al.<sup>1,2</sup>; Zhuang, Hui<sup>3</sup>; Danner, B.; Kunkin, L.<sup>1,2</sup>; Sachet, A.; Shukla, R.; Berg, Sonn<sup>4</sup>; David J. Bresnahan<sup>5</sup>; Wandi Zhang<sup>6</sup>; Adrienne Luoma<sup>7</sup>; Andrea Gobbi<sup>8</sup>; Dennis B. Kunkin<sup>1,2</sup>; Steven Peter<sup>9</sup>; Christine Ober<sup>10</sup>; Olaf Oltmer<sup>10</sup>; David C. Carter<sup>11</sup>; Shuang Li<sup>12</sup>; David J. Lieb<sup>13</sup>; Thomas Eisenhauer<sup>14</sup>; Eva Gijsel<sup>15</sup>; Jonathan Seaman<sup>16</sup>; Elizabeth L. Linsen<sup>17</sup>; Juval Juvel<sup>18</sup>; Källbergapannen Nellapragash<sup>19</sup>; Andrew M. Salazar<sup>20</sup>; Heather Daley<sup>21</sup>; Michael Stephan<sup>22</sup>; Elizabeth L. Bichakjian<sup>23</sup>; Charles H. Eberle<sup>24</sup>; Mark A. Hammon<sup>25</sup>; Nadia Neves<sup>26</sup>; Steven Ritter<sup>27</sup>; Scott J. Sander<sup>28</sup>; Daniel H. Barnard<sup>29,30</sup>; Jon C. Aster<sup>31,32</sup>; B. Catherine I. Wu<sup>33</sup>; Kai Weng<sup>34</sup>; Bernadette Neuberg<sup>35</sup>; Dennis Ritter<sup>36</sup>; Steven Ritter<sup>37</sup>; Eric S. Landig<sup>38</sup>; Edward F. Pratichetti<sup>39</sup>; Nuri Haouzi<sup>40</sup>

“All indels were manually reviewed in Integrative Genomics Viewer.”

“Coverage of at least 50 reads in both tumor and normal samples, >20% of reads supporting the variant in tumor samples, and <5% of reads supporting the variant in normal samples.”



## RNF43 is frequently mutated in colorectal and endometrial cancers

Marios Giannakakis<sup>1,2,3</sup>, Eren Hudai<sup>1,3,4,5</sup>, Ximeng Jieming Mu<sup>1,2</sup>, Mai Tamashii<sup>1</sup>, Joseph Rosenblatt<sup>1</sup>, Kristian Cibulskis<sup>1</sup>, Gordon Jaksica<sup>1</sup>, Michael S Lawrence<sup>1</sup>, Zhi Ran Quan<sup>1</sup>, Reiko Nishihara<sup>1,2,6</sup>, Elsner M Van Allen<sup>1,2,7</sup>, William C Hahn<sup>1,2</sup>, Stacy B Gottsch<sup>1</sup>, Eric S Lander<sup>1,2,8</sup>, Gad Getz<sup>1,2,9</sup>, Shaili Ogilvie<sup>1,2,10</sup>, Charles Fuchs<sup>1,2,11</sup> & Levi A Garraway<sup>1,2,12</sup>

collected from participants in 2 prospective cohort studies, the Nurses' Health Study (NHS) and the Health Professionals Follow-up Study (HPFS),<sup>1</sup> showed mutations in 10% (18/182) cases of breast cancer in 23% of tumors (0.1–0.6%).<sup>2</sup> Supplementary Table 11. Insertions or deletions involving p53Arg175Arg constituting frameshift mutations resulting in p53Arg175Arg (missense) instability (MSI) loci of seven and six cGPs, respectively, accounted for 31.7% ( $n=6$ ) and 8.3% ( $n=1$ ). Arg175 of the ENK2-43 mutations identified [Fig. 1A]. To exclude the possibility that these mutations represented technical artifacts, we validated 31 of the ENK2-43 mutations (37% of the 82 reactions) that had leftover DNA available and achieved coverage of  $\geq 80\%$  in

- Rarely are the methods/criteria for manual review described
  - Definitely not standardized

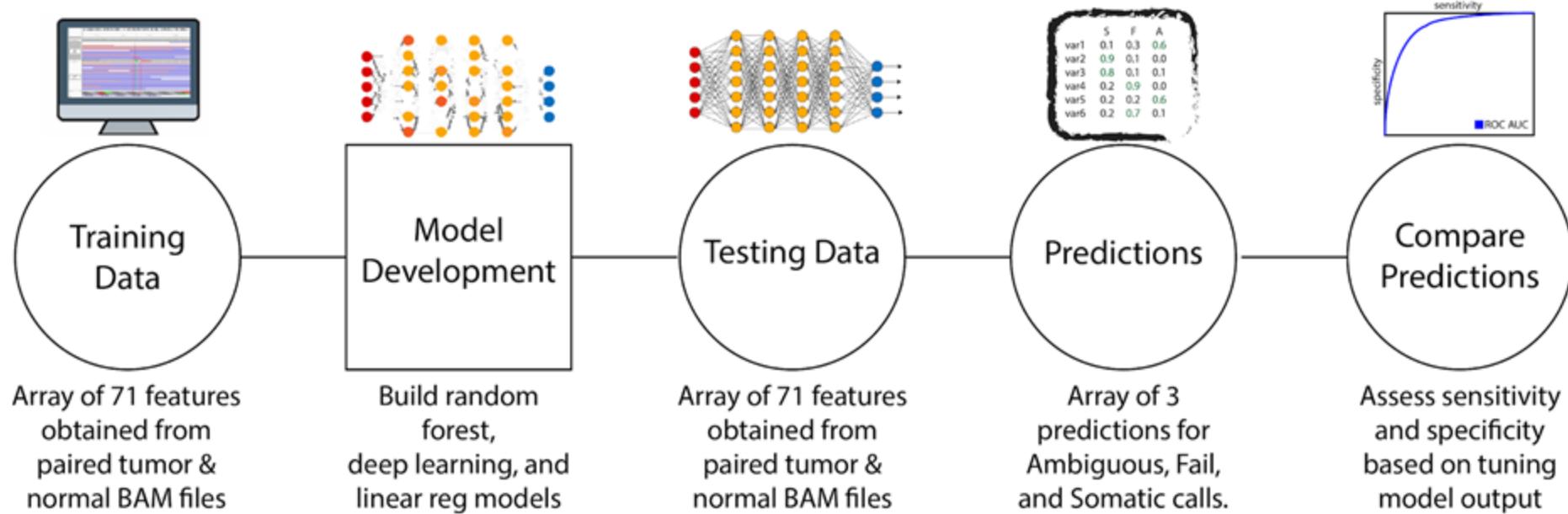


## Manual review is a largely undocumented source of error, suffers from inter-reviewer variability, and requires significant resources

- We completed a study of independent/blinded manual review of a set of variants.
- High rates of inter-reviewer variability
  - Only "fair agreement" - Kappa = 0.37
  - 77.3% good (3/3 reviewers) or acceptable (2/3 reviewers) agreement on variant calls
- Manual review is labor-intensive, time-consuming
  - ~70-100 variants per hour



# Can manual review be automated with AI?



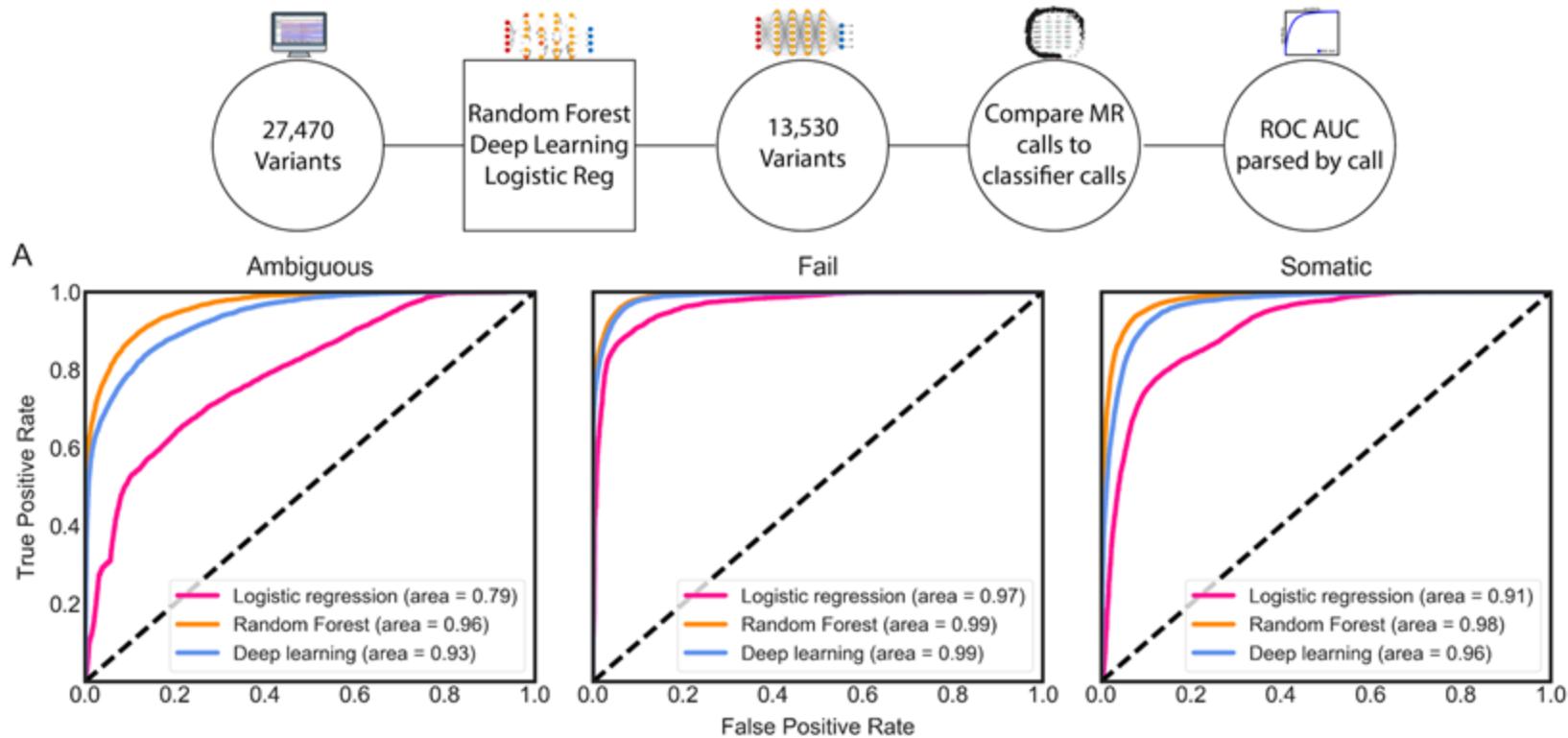
# Assembled database of 40,000+ manually reviewed variants

	Training Set (n=27,470)	Testing Set (n=13,530)	Total (n=41,000)
<b>Malignancy (410 cases)</b>			
Leukemia (243 cases)	5,815	2,877	8,692
Lymphoma (23 cases)	1,263	628	1,891
Breast (135 cases)	8,986	4,320	13,306
Small Cell Lung (18 cases)	9,177	4,601	13,778
Glioblastoma (17 cases)	844	412	1,256
Melanoma (1 cases)	185	100	285
Colorectal (1 case)	842	419	1,261
Gastrointestinal Stromal (1 case)	70	31	101
Malignant Peripheral Nerve Sheath (1 case)	288	142	430
<b>Sequencing Methods</b>			
Capture Sequencing	9,479	4,755	14,234
Exome Sequencing	9,367	4,677	14,044
Whole Genome Sequencing	8,624	4,098	12,722
<b>Variant Calls</b>			
Somatic	12,266	6,115	18,381
Ambiguous	7,189	3,454	10,643
Fail	5,909	2,945	8,854
Germline	2,106	1,016	3,122

- Nine different tumor subtypes tested
  - 266 hematologic tumors and 144 solid tumors
- Multiple sequencing platforms tested
- Even distribution of variant calls (Somatic, Germline, Ambiguous, or Fail)
- Features used by manual reviewers extracted from alignments



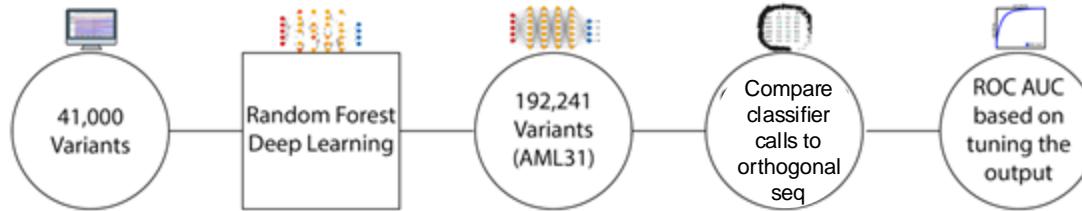
# Highly accurate (ROC AUC > 0.9) predictive model



- **Random Forest and Deep Learning models out-perform Logistic Regression model**
- **Ambiguous calls are the most difficult to classify**



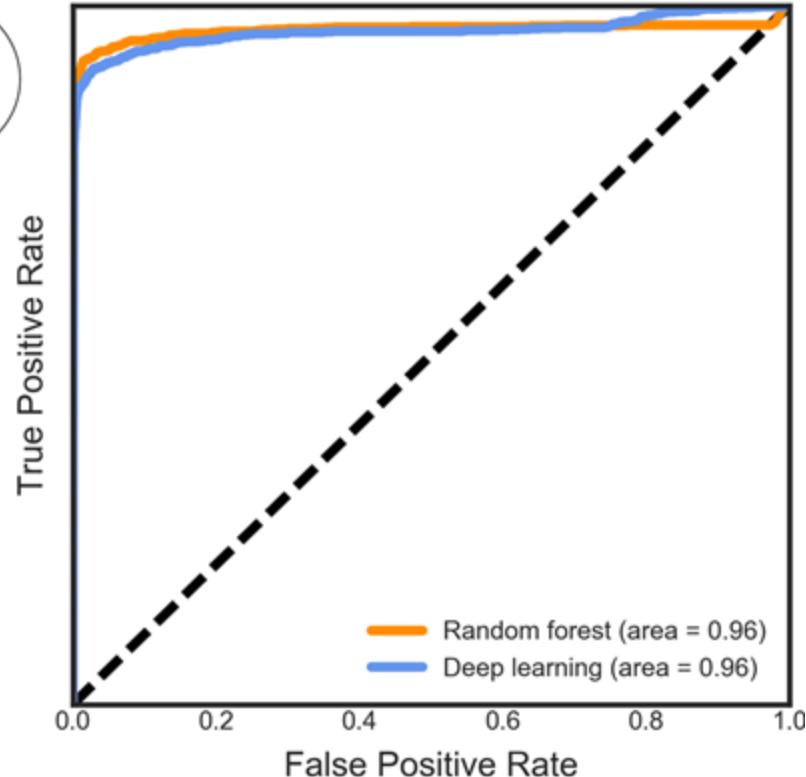
# Excellent performance on Orthogonally validated variants



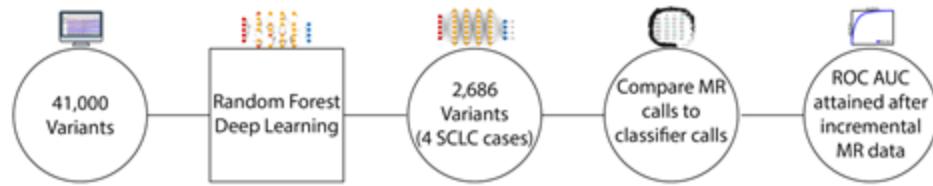
## Orthogonal Validation Strategy:

- 192,241 variants were called by automated somatic variant callers.
- 1,343 true positives and 190,898 false positives were identified by 1000X sequencing
- Model predictions were compared to orthogonal sequencing results

**Use of orthogonal sequencing data to classify variants performed exceptionally well**



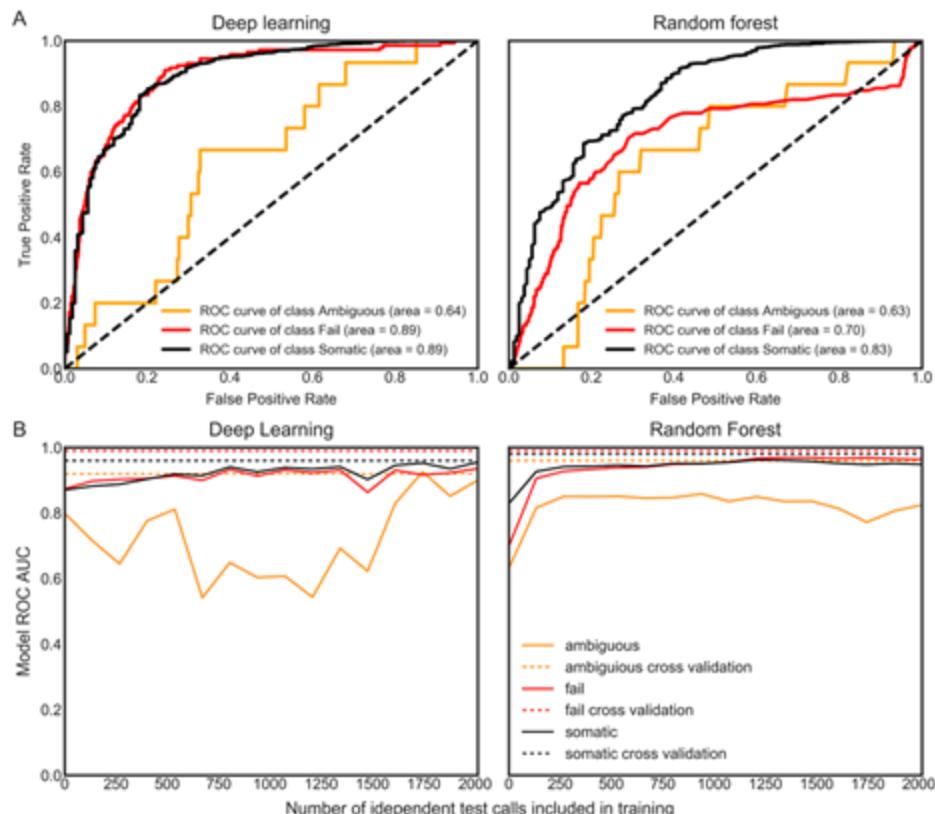
# Performance is reduced on external datasets but can be recovered with a modest amount of re-training



“External Lab” Validation Strategy:

- Analyze performance on 4 SCLC cases with different manual reviewers, different sequencing platform, and different automated variant callers.
- Initial AUC showed suboptimal performance.
- Model was re-trained with 5% of manual review calls from new sequencing data.

Retraining the data with ~200 variants shows exceptional recovery of model performance.



# A deep learning approach to automate refinement of somatic variant calling from cancer sequencing data



Josh  
Swamidass

Benjamin J. Ainscough  <sup>1,2,12</sup>, Erica K. Barnell  <sup>1,12</sup>, Peter Ronning<sup>1</sup>, Katie M. Campbell  <sup>1</sup>, Alex H. Wagner  <sup>1</sup>, Todd A. Fehniger  <sup>2,3</sup>, Gavin P. Dunn<sup>4</sup>, Ravindra Uppaluri<sup>5</sup>, Ramaswamy Govindan<sup>2,3</sup>, Thomas E. Rohan<sup>6</sup>, Malachi Griffith  <sup>1,2,3,7</sup>, Elaine R. Mardis<sup>8,9</sup>, S. Joshua Swamidass<sup>10,11\*</sup> and Obi L. Griffith  <sup>1,2,3,7\*</sup>



Benjamin  
Ainscough

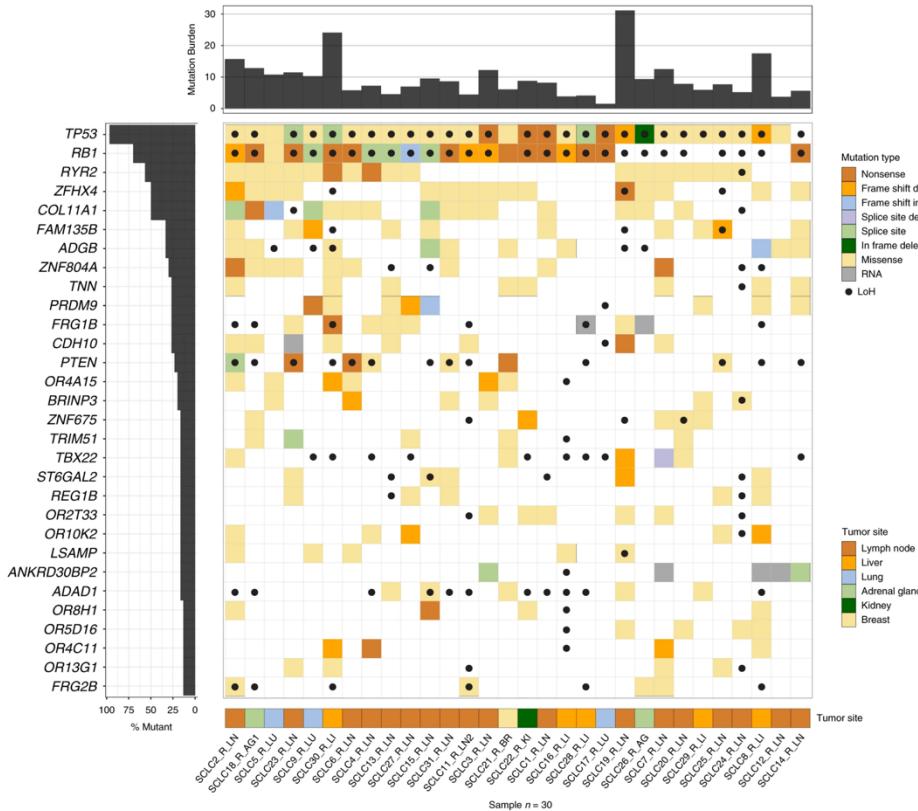
Cancer genomic analysis requires accurate identification of somatic variants in sequencing data. Manual review to refine somatic variant calls is required as a final step after automated processing. However, manual variant refinement is time-consuming, costly, poorly standardized, and non-reproducible. Here, we systematized and standardized somatic variant refinement using a machine learning approach. The final model incorporates 41,000 variants from 440 sequencing cases. This model accurately recapitulated manual refinement labels for three independent testing sets (13,579 variants) and accurately predicted somatic variants confirmed by orthogonal validation sequencing data (212,158 variants). The model improves on manual somatic refinement by reducing bias on calls otherwise subject to high inter-reviewer variability.



Erica  
Barnell



# We have applied this model to multiple large-scale cancer genome projects, including our survey of relapsed SCLC



- High mutation burden requiring review of 100s or 1000s of variants
- Confirmed TP53 and RB1 nearly universally co-altered in SCLC
- COL11A1 next most significantly mutated, may mediate resistance to platinum chemotherapy
- Relapse samples also characterized by recurrent alterations in WNT signaling



ARTICLE

DOI: 10.1038/s41467-018-06162-9

OPEN

# Recurrent WNT pathway alterations are frequent in relapsed small cell lung cancer

Alex H. Wagner<sup>1</sup>, Siddhartha Devarakonda<sup>2,3</sup>, Zachary L. Skidmore<sup>1</sup>, Kilannin Krysiak<sup>1,2</sup>, Avinash Ramu<sup>1</sup>, Lee Trani<sup>1</sup>, Jason Kunisaki<sup>1</sup>, Ashiq Masood<sup>2,3,9</sup>, Saiama N. Waqar<sup>2,3</sup>, Nicholas C. Spies<sup>1</sup>, Daniel Morgensztern<sup>2,3</sup>, Jason Waligorski<sup>1</sup>, Jennifer Ponce<sup>1</sup>, Robert S. Fulton<sup>1</sup>, Leonard B. Maggi Jr.<sup>2,3,4</sup>, Jason D. Weber<sup>2,3,4</sup>, Mark A. Watson<sup>3</sup>, Christopher J. O'Conor<sup>5</sup>, Jon H. Ritter<sup>5</sup>, Rachelle R. Olsen<sup>6</sup>, Haixia Cheng<sup>6</sup>, Anandaroop Mukhopadhyay<sup>6</sup>, Ismail Can<sup>1</sup>, Melissa H. Cessna<sup>7</sup>, Trudy G. Oliver<sup>6</sup>, Elaine R. Mardis<sup>1,3,8,10</sup>, Richard K. Wilson<sup>1,3,8,10</sup>, Malachi Griffith<sup>1,2,3,8</sup>, Obi L. Griffith<sup>1,2,3,8</sup> & Ramaswamy Govindan<sup>2,3</sup>



Alex  
Wagner



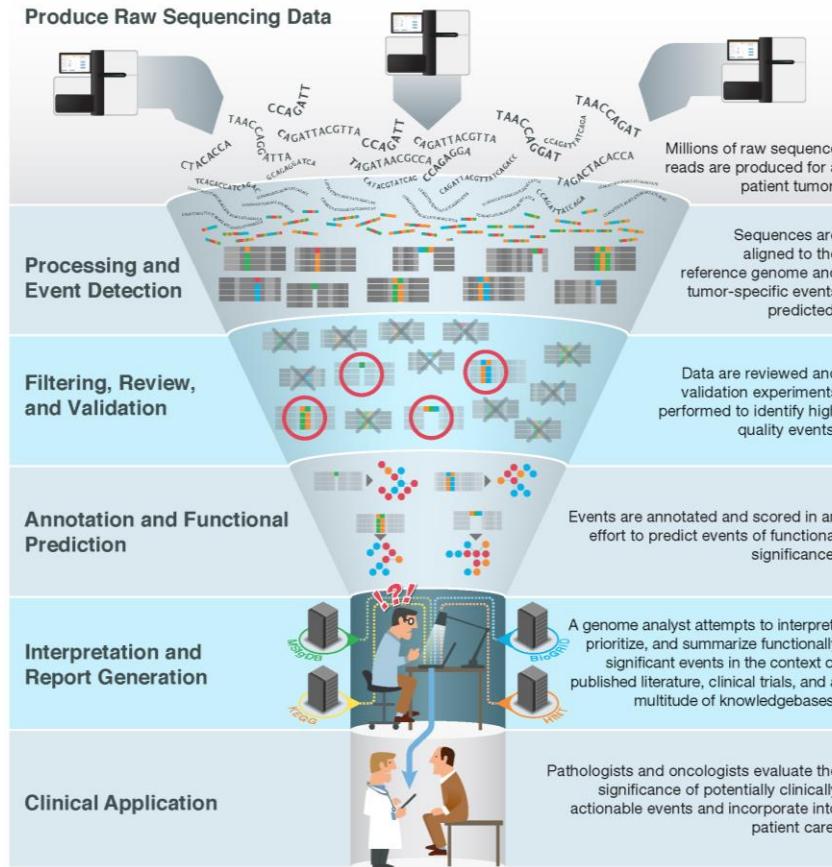
Ramaswamy  
Govindan



Siddhartha  
Devarakonda



# How else can AI be applied to genomic medicine?



Variant calling

Manual review

Functional prediction

Interpretation

Decision-making



## **II) Interpretation - how do we know if a cancer variant is clinically significant?**



# What is a variant of \*significance\* in cancer?

- **Predictive** of therapeutic response
  - *BRAF V600E predicts sensitivity to vemurafenib*
- **Diagnostic** of tumor subtype
  - *DNAJB1-PRKACA fusion differentiates fibrolamellar hepatocellular carcinoma from conventional HCC*
- **Prognostic** of survival change
  - *TP53 mutations are associated with worse progression-free survival in lung adenocarcinoma*
- **Predisposing** for cancer development
  - *Patients with the RUNX1 Y260\* mutation are associated with increased risk of developing acute myeloid leukemia*



# CIViC: an open knowledgebase and curation system for clinical interpretation of variants in cancer

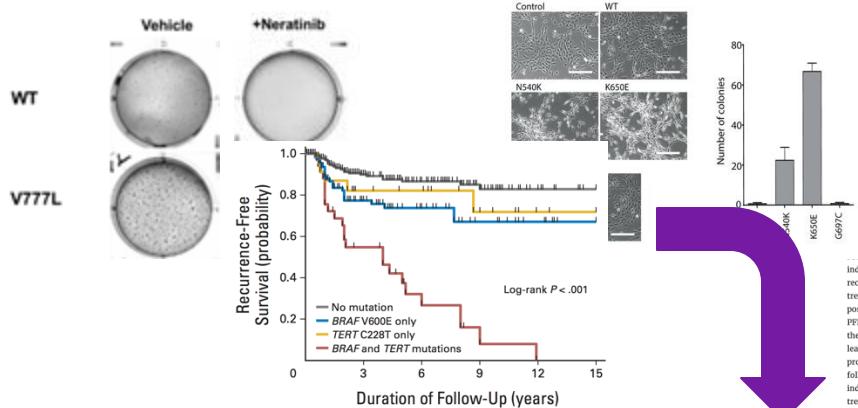


[www.civicdb.org](http://www.civicdb.org)

Griffith et al. Nat. Gen. 2017

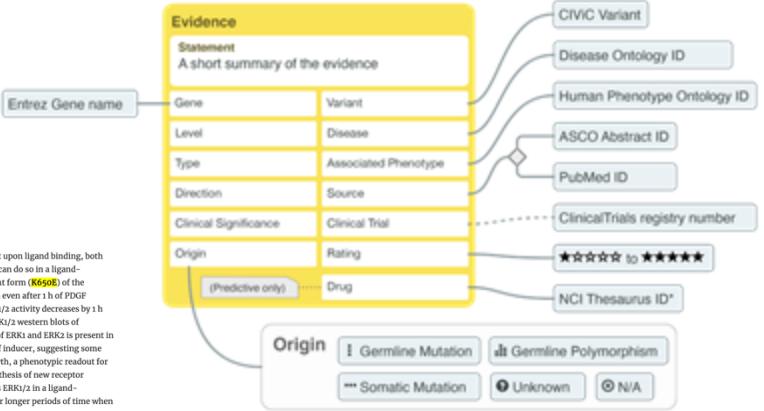


# Heterogeneous evidence



->

# Structured data model



**CIViC** CLINICAL INTERPRETATION OF VARIANTS IN CANCER

Quicksearch

Evidence / EID288 / Summary

EID288

Summary Comments Revisions Flags Events

Description

In MCF10A cell lines, the V777L mutation was shown to be sensitive to neratinib.

Status Accepted Submitted Jun 21, 2015 by NickSpies Accepted Jun 21, 2015 by kkrysiak

Source PubMed: Bose et al., 2013, Cancer Discov

Clinical Trial None Specified

Type Predictive Direction Supports

Clinical Significance Sensitivity / Response Disease Phenotype

Level Rating ★★★★

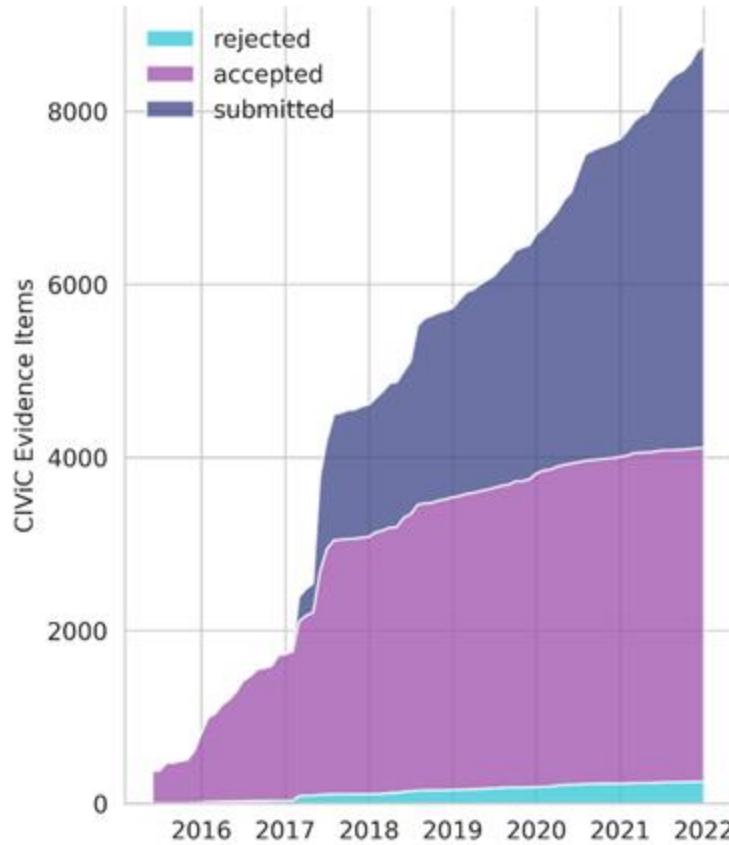
Breast Cancer None Specified

Drug Neratinib

651 ObiGriffith

Curators: Editors:

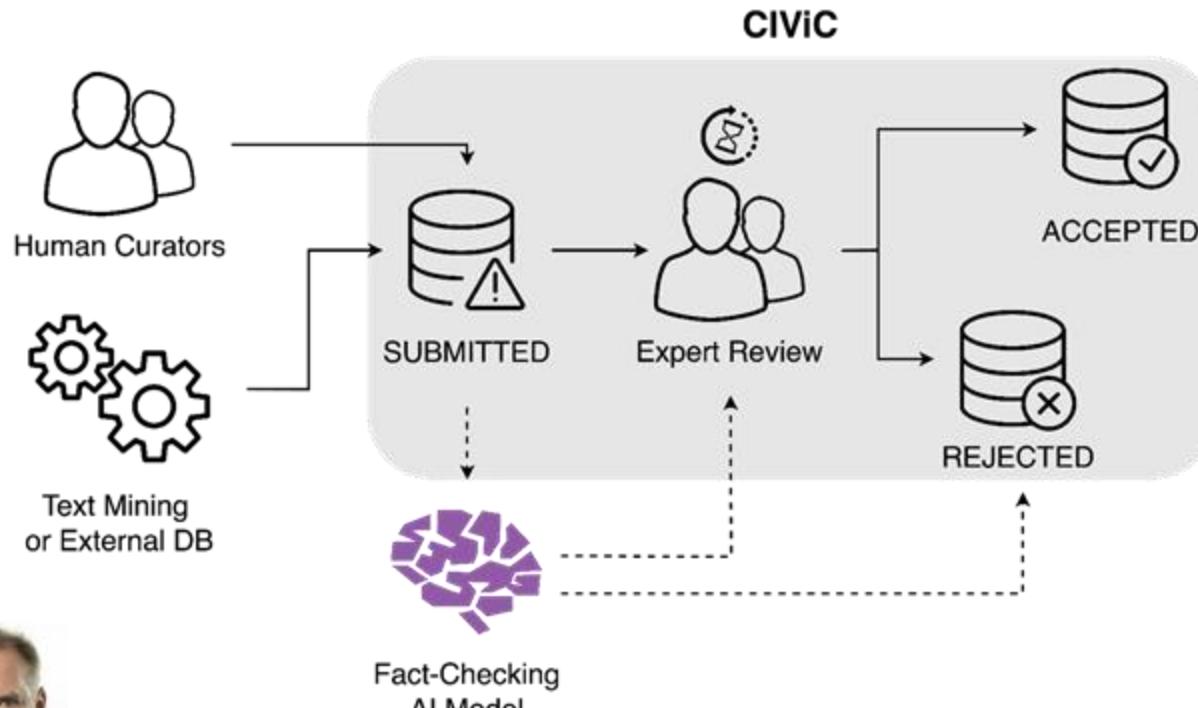
# Expert moderation of submitted content is a major bottleneck



- Entries are being submitted to CIViC faster than experts can review them (right)
- Automated fact-checking could alleviate some of this burden
- Dearth of open-data scientific datasets in fact-checking and natural language interpretation



# Proposed Integration of Fact Checking AI into CIViC



Caralyn Reisle (Steven Jones): Poster 64

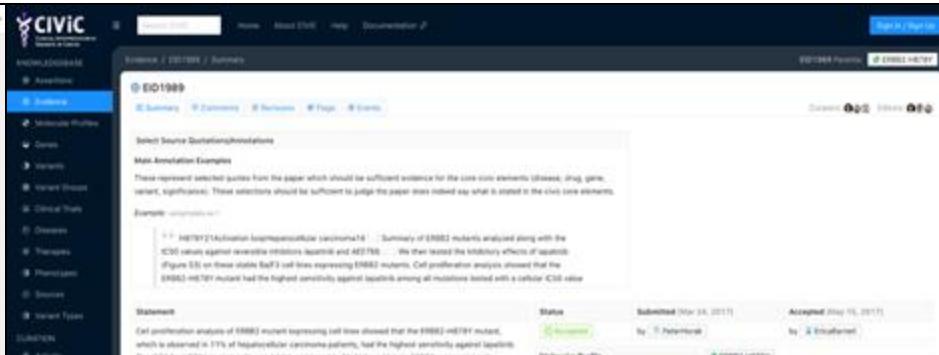


# Annotation Process: Creating the Cancer-Facts Dataset

downstream effectors in wild-type or mutated ALK-expressing Ba/F3 cells depleted of IL-3 for 6 hours. The mobilities of molecular weight (MW) standards are shown on the left. *a*, Growth of mutated ALK-expressing Ba/F3 cells exposed to TAE684 for 72 hours. The values are mean  $\pm$  SD of triplicate experiments.

The Ba/F3 assay has been validated for a broad spectrum of oncogenic tyrosine kinase alleles including mutant EGFR<sup>12</sup> and FLT3<sup>10</sup>, and thus we treated Ba/F3 cells expressing each of the ALK mutations with increasing concentrations of TAE684, a highly potent ALK inhibitor<sup>11,13</sup>. The activating mutation, F1174L, was found to be extremely sensitive to TAE684, with an IC<sub>50</sub> of 8 nM, identical to that of NPM-ALK-expressing Ba/F3 cells. The R1275Q mutation was also sensitive to TAE684, albeit with a much higher IC<sub>50</sub> of 328 nM. In contrast, Ba/F3 cells expressing FLT3-ITD or wild-type ALK, did not respond to TAE684 (IC<sub>50</sub> 4.5  $\mu$ M; Fig. 1*a*).

Analysis of the ALK gene in a panel of 30 neuroblastoma cell lines revealed sequence variants in 6, including 3 different cell lines containing the F1174L mutation (Kelly, SH-SY5Y and LAN-1), which was also the most common mutation in the primary tumors (Table 1). An R1275Q mutation, identical to the one found in primary sample 411, was also detected in the SMS-KCNR cell line. We observed dose-dependent growth inhibition of the SH-SY5Y (F1174L) and Kelly (F1174L) neuroblastoma cell lines with increasing concentrations of TAE684, (IC<sub>50</sub> of 258 and 416 nM respectively; Fig. 1*b*). These results are



## Evidence is selected using hypothes.is



Hypothes.is Selections are Displayed in CIViC via a Custom Chrome Extension

Claim/Evidence Pair Labels are Confirmed with a 2nd Round of Annotation



# Evidence Item Consistency Analysis

Significance	<input checked="" type="checkbox"/> Sensitivity / Response
Type	<input checked="" type="checkbox"/> Predictive
Molecular Profile	<input checked="" type="checkbox"/> MGMT Promoter Methylation
Disease	<input checked="" type="checkbox"/> Glioblastoma
Phenotype	Not specified
Therapy	<input checked="" type="checkbox"/> Temozolomide

Statement

In a randomized clinical trial, patients with MGMT promoter methylation benefitted from temozolomide. This benefit was also methylation status dependent, as those without methylation did not see increased survival.

Direction

Supports

## Prompt

Given a claim and some evidence, determine the stance of the evidence toward the claim. Use the following labels: SUPPORTS, REFUTES, and NEI (not enough information). Give your reasoning.

claim: molecular profile is <MP>; evidence type is <evidence type>; significance is <significance>; therapy is <therapy>; disease is <disease>

## Expected Response

evidence: <statement>

label: <direction>



# Example #1: AI fact-checking model flags errors in evidence direction field



EID7393

Summary Comments Revisions Flags Events

Status Submitted (Jul 8, 2019) Accepted (Jul 9, 2019)

Molecular Profile KIT V560\_L576DEL

MP Expression KIT V560\_L576DEL

Disease Gastrointestinal Stromal Tumor

Phenotype Not specified

Therapies Ponatinib Imatinib

Therapy Interaction Type Substitutes

**Statement**

In the patient-derived cell line GIST430 harboring KIT V560\_L576del primary activating mutation implanted in mice, imatinib and ponatinib were shown to reduce KIT phosphorylation compared to the vehicle. In GIST430 cell lines, ponatinib and imatinib reduced phospho-AKT and phospho-ERK levels. A cell line derived from GIST430 patients expressing KIT V560\_L576del primary mutation demonstrated sensitivity to imatinib ( $IC_{50}$ : 6.1nmol/L) and ponatinib ( $IC_{50}$ : 12nmol/L) treatments. In comparison, the insensitive KIT-independent GIST226 cell line showed  $IC_{50}$  >5000nmol/L and 2807nmol/L for these drugs respectively.  $IC_{50}$  was determined by assessing cell viability.

Type Predictive Direction Does Not Support

Significance Sensitivity / Response Variant Origin Somatic

Level Preclinical evidence Rating ★★★★☆

[mosaicml\_mpt-7b-chat]: **SUPPORTS**. The evidence suggests that both imatinib and ponatinib effectively target GIST molecular profile with high sensitivity . . .

[BioMistral\_BioMistral-7B-DARE]: **SUPPORTS**. The cell line derived from GIST430 patients expressing KIT V560\_L576del primary mutation demonstrated sensitivity to imatinib and ponatinib treatments . . .



RID60338 Evidence Direction Updated By kkrysiak for The McDonnell G

Original Evidence Direction Evidence Direction Diff

Does Not Support Supports



## Example #2: AI fact-checking model flags errors in significance field

EID1516

Summary Comments Revisions Flags Events Curators: 

**Statement**  
Patients with a triple mutation of NPM1, DNMT3A, and NRAS (G12/13) (three-way interaction p=0.04) showed a relatively benign prognosis with an 8-year survival rate of 75% and a 10-year survival rate of 55%. Conversely, patients with a concomitant NPM1/DNMT3A co-mutation, in the absence of an NRAS (G12/13) mutation, showed a clinically unfavorable 10-year survival rate of 30%.

Type	<input checked="" type="checkbox"/> Prognostic	Direction	<input checked="" type="checkbox"/> Supports	Status	Submitted (Jul 7, 2016)	Accepted (Jul 7, 2)
Significance	<input checked="" type="checkbox"/> Poor Outcome	Variant Origin	<input checked="" type="checkbox"/> Somatic	by  EricaBarnell	by  NickSpies	
Level	B - Clinical evidence	Rating		Molecular Profile	 NRAS G12/G13	
Source	PubMed: Papaemmanull et al., 2016			MP Expression	 NRAS  G12/G13	
Clinical Trial	 NCT00146120			Disease	 Acute Myeloid Leukemia	
				Phenotype	Not specified	
				Therapy	Not applicable	



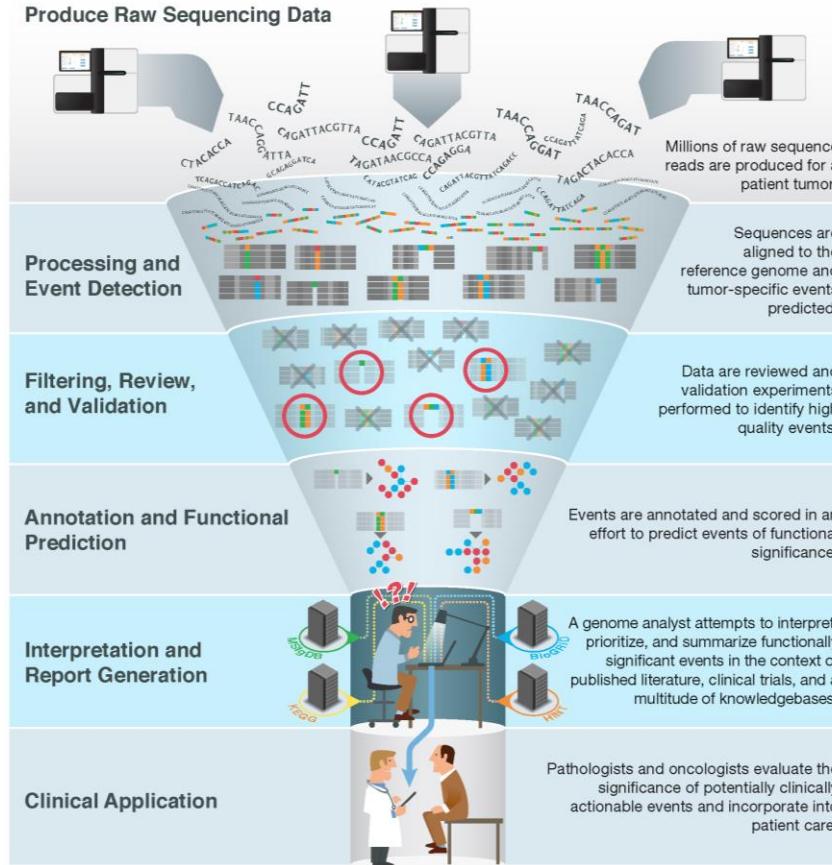
[BioMistral\_\_BioMistral-7B-DARE]: All elements are present. Patients with NRAS G12/13 mutation had better outcome (8 yr OS 75%) than those without this mutation (10 yr OS 30%) . . .



# How else can AI be applied to genomic medicine?

Communication

- figures
- grants
- letters
- homework



Variant calling

Manual review

Functional prediction

Interpretation

Decision-making



# What can LLMs help with?

The screenshot shows a sidebar with a list of previous interactions:

- Yesterday
- Missouri Doc Fees Kia EV9
- St Louis County Vehicle Tax
- How old am I

Below this, there are sections for "Previous 7 Days" and "Previous 30 Days", each listing several topics.

At the bottom of the sidebar, there is a "Temporary" button and a "PLUS" user icon.

The main interface features a search bar with the placeholder "Ask anything". Below the search bar are buttons for "+", "Search", "Deep research", and "...". To the right is a microphone icon for voice input.

**What can I help with?**

Why am I trying to do this without using LLMs?

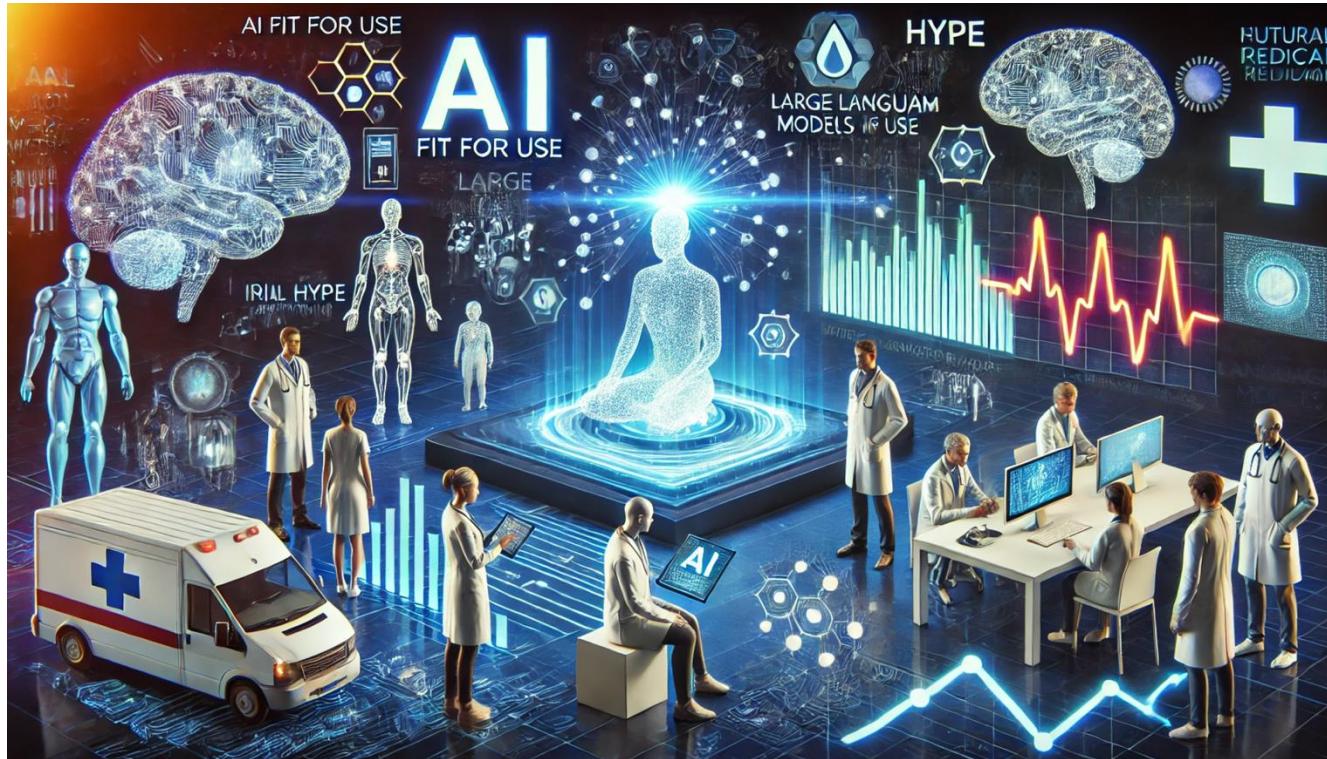


# I asked DALL-E to make a figure for a review paper

I would like to create a high-level conceptual image to summarize a paper I wrote. The paper is titled "Title: Looking Forward to AI and Medicine: Where Are We, and Where Are We Going?" The key themes of the paper are: (1) Why this moment of AI hype could be different. The emergence of LLMs; (2) The importance of Fit for use - why this should be evaluated differently; (3) The need for evaluation - how it is different in medicine (data sharing, regulation, etc); (4) What we expect to happen - After initial hype, we enter a period of diminishing returns in the short/medium term but then improvements accelerate eventually leading to transformation of many areas of medical practice.



# Interesting but weird - what are these people doing?



Please create new versions without human figures. Something more conceptual



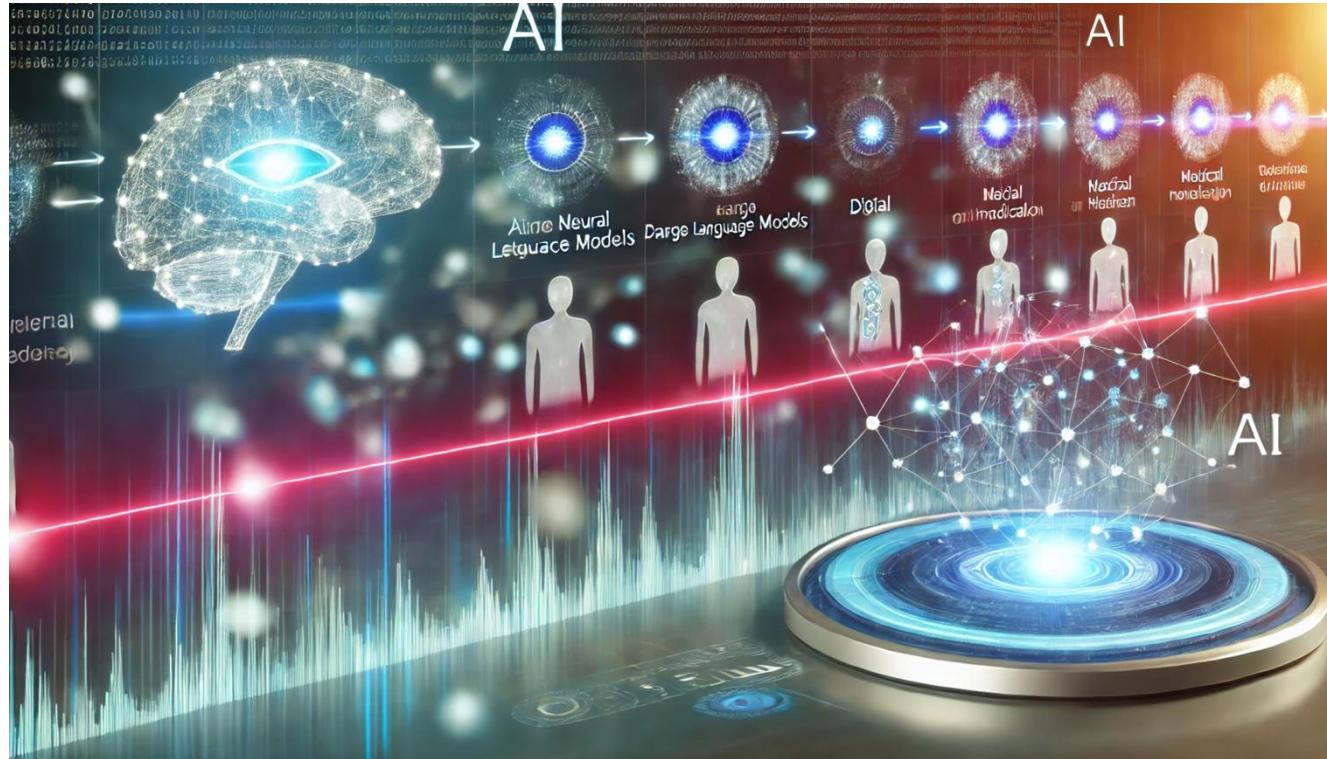
# Cooler but with some problem elements



This is great. However, I don't like the red line going down the person on the right. Also, please remove the word "Hype" from the image

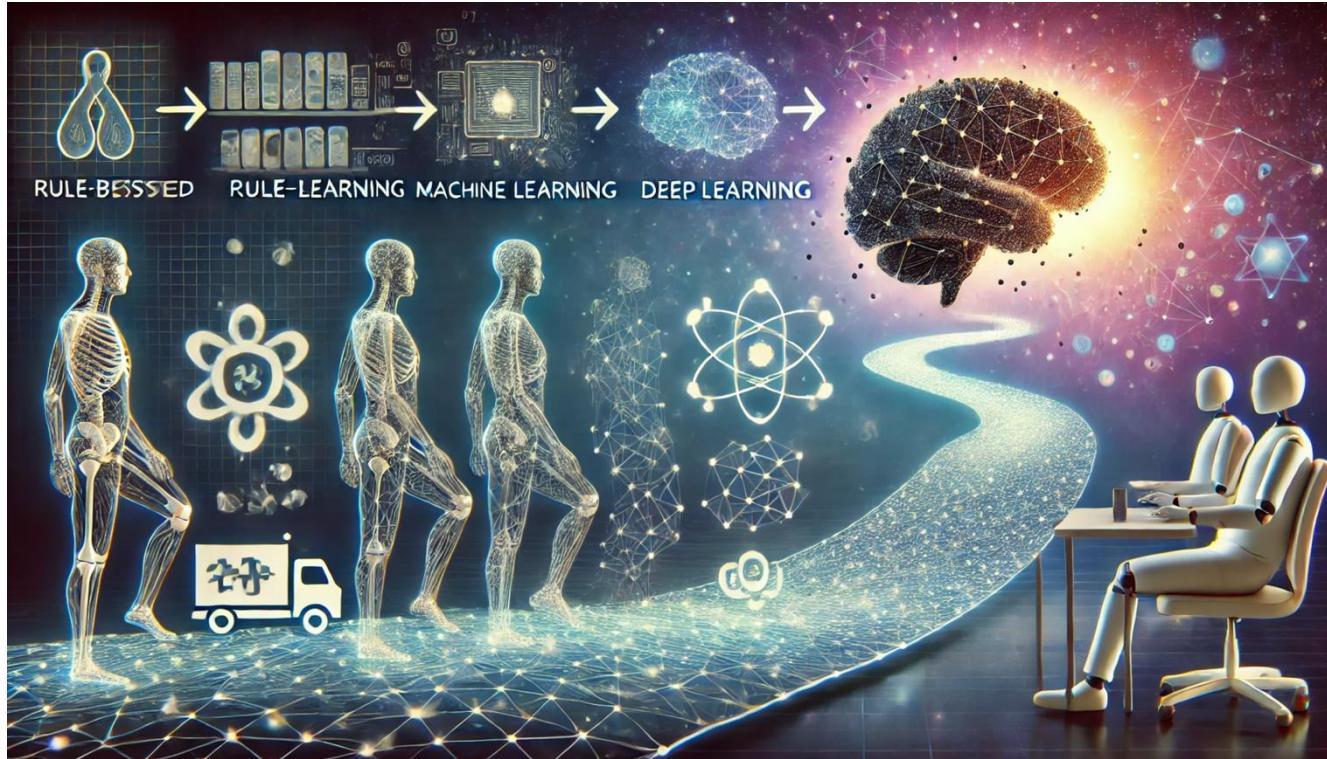


# Very different given minor suggestion, but interesting new direction



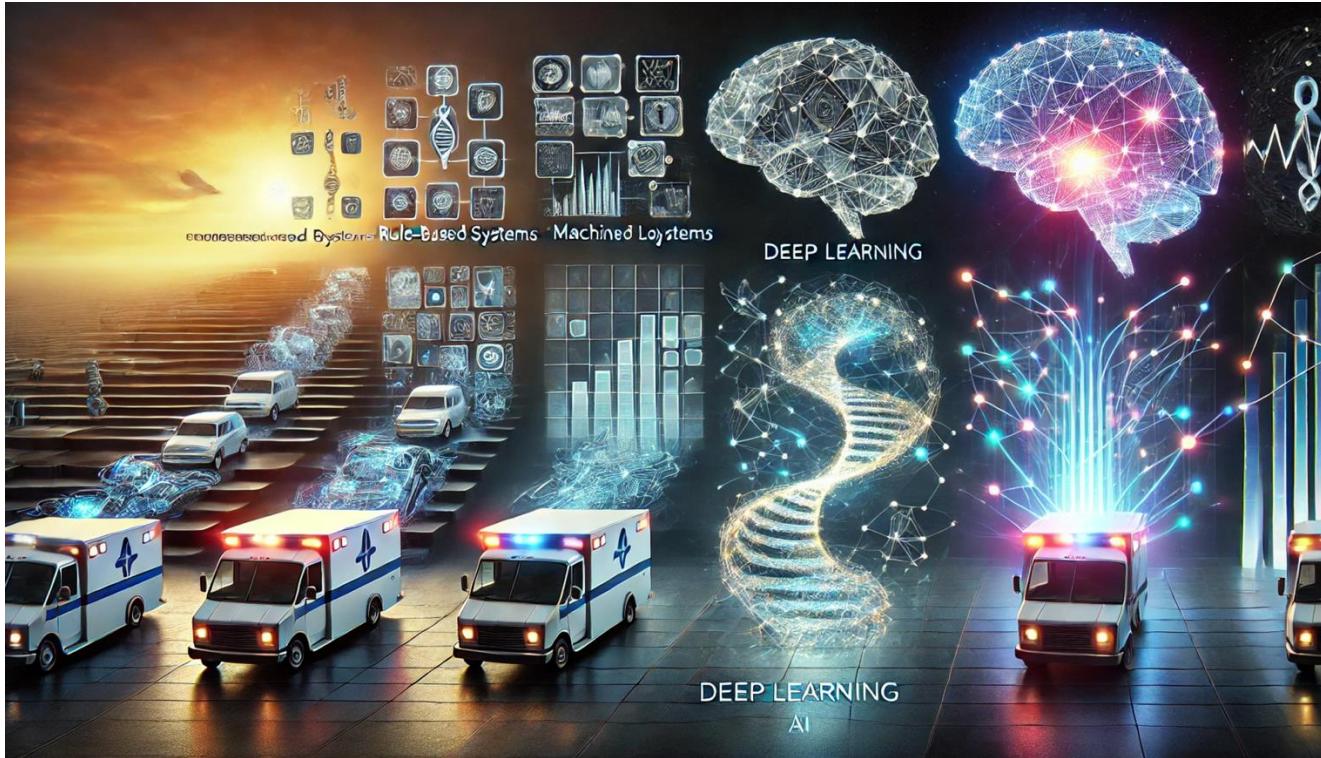
I like this last version which includes both AI concepts and a timeline. Create similar figures which graphically depict the evolution from rule-based systems -> machine learning -> deep learning -> generative AI





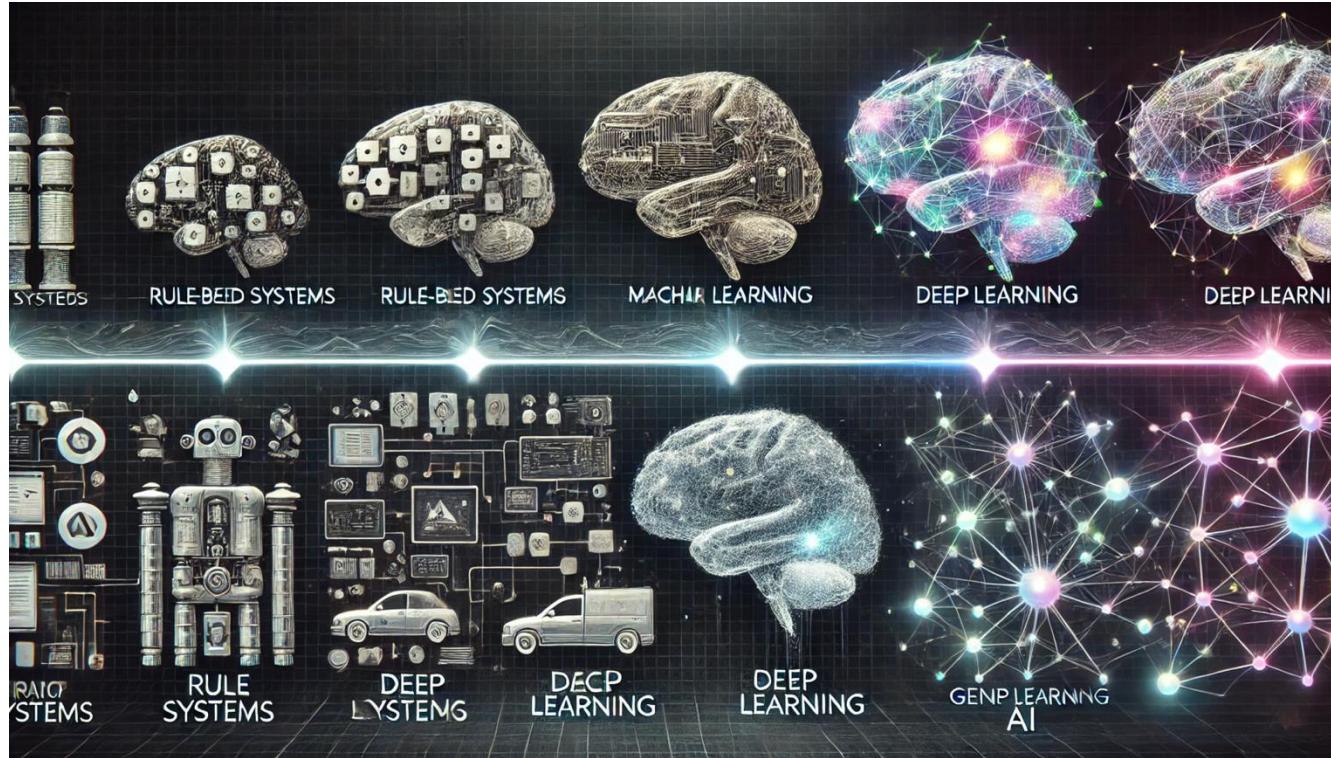
These are OK but please don't include any vehicles, include conceptual elements more like the previous two images, and make sure that the four types of AI are correctly spelled and represented in order: (1) "rule-based systems"; (2) "machine learning"; (3) "deep learning"; (4) "generative AI"





These are better but now create versions without vehicles and just limit to the 4 stages of AI

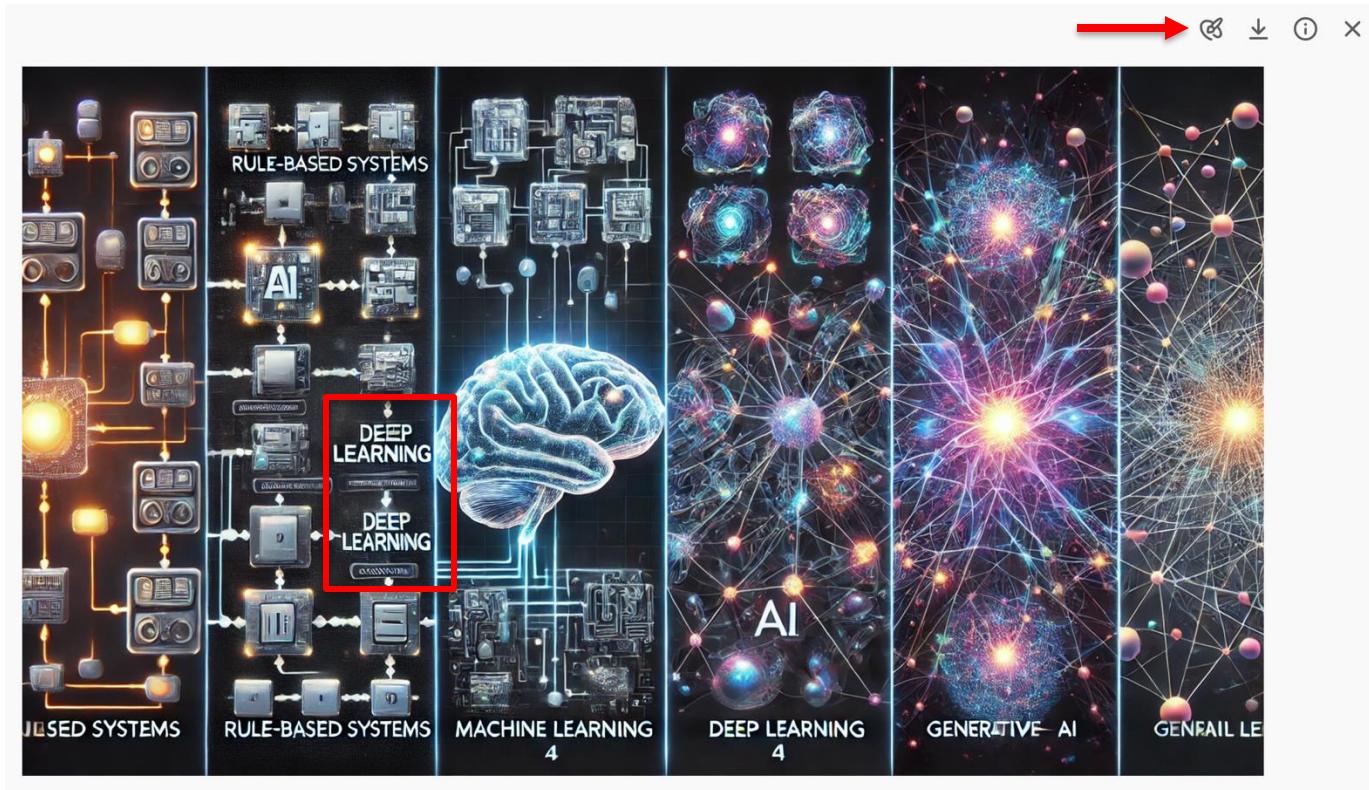




These look great but still don't have the 4 stages, correctly spelled, just appearing once each. Please just depict these four things in order: (1) "rule-based systems"; (2) "machine learning"; (3) "deep learning"; (4) "generative AI"



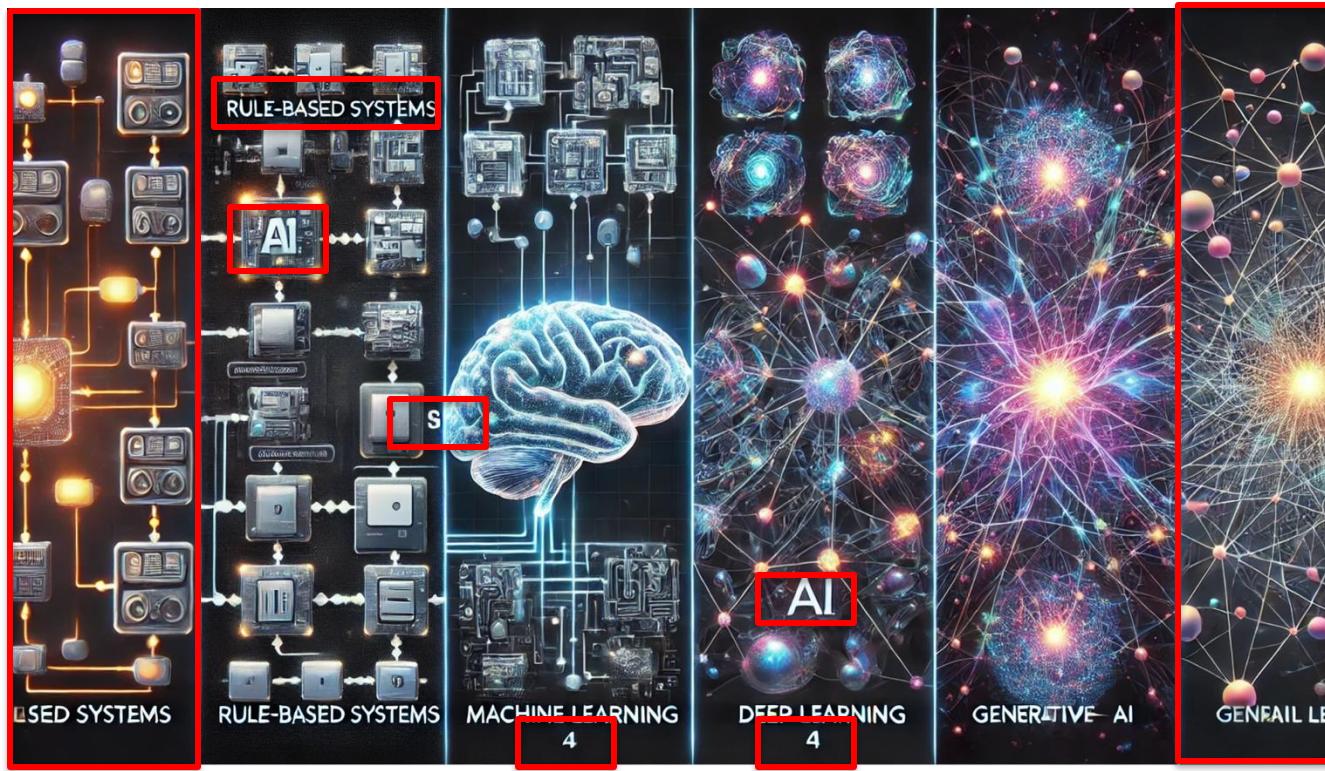
# "Select and area to edit"



remove the words deep learning from this section and replace with graphics similar to other parts of the panel



~10 more iterations



Remove the, outer panels and highlighted letters but otherwise leave the image as is



# Final Figure: The different stages of artificial intelligence

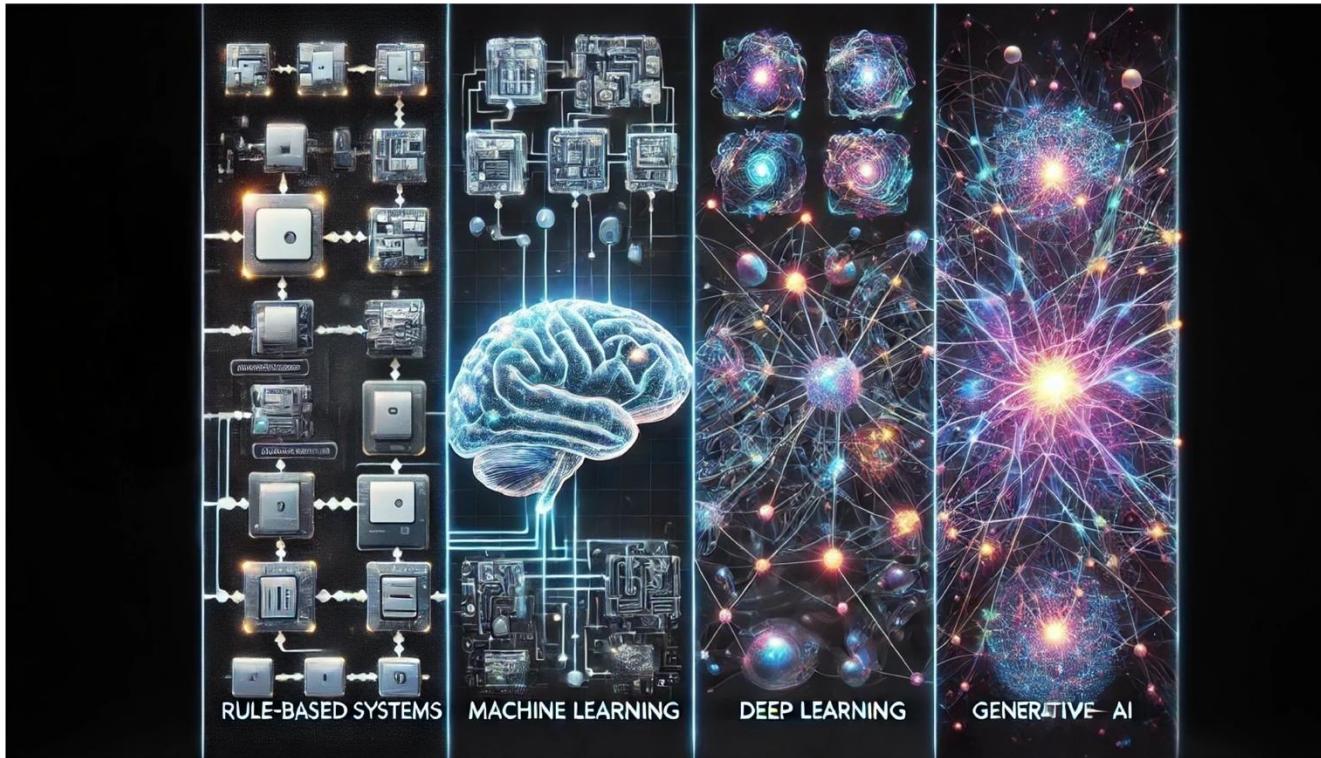


Figure 1. The different stages of artificial intelligence. This image was created (with significant prompt iteration and targeted edit requests) by DALL-E, an AI image generator.



# Published in Missouri Medicine - Jan/Feb 2025

SCIENCE OF MEDICINE | FEATURE SERIES

## Looking Forward to AI and Medicine: Where Are We, and Where Are We Going?

by Adam Wilcox, PhD, Malachi Griffith, PhD & Obi Griffith, PhD



We speculate that the integration of AI, particularly generative AI, into medicine is expected to progress at a slower pace than in fields such as research.



Adam Wilcox, PhD, FACP, (pictured), is Director, Center for Applied Health Informatics, Professor of Medicine, Division of General Medicine & Geriatrics, Washington University School of Medicine, St. Louis, Missouri. Malachi Griffith, PhD, and Obi Griffith, PhD, are both Associate Professors of Medicine (Oncology) and Genetics and Assistant Director of the McDonnell Genome Institute at Washington University, St. Louis, Missouri.

### Abstract

Artificial intelligence (AI) has emerged as a significant area of interest in medicine, with the potential to influence various aspects of health care. However, the real benefits of applications of AI to medicine often become obscured by the considerable attention, and hype, around its capabilities. To better understand AI's role in medicine, it is important to contextualize its development and review the stages of AI evolution that have contributed to its present state.

### Introduction

My (AW) first exposure to AI was taking a course in the topic as part of my graduate studies in a new field called "Medical Informatics." Most of the students in the class were training as computer scientists, and during that semester were particularly interested in a chess match between then-world chess champion Garry Kasparov and IBM's Deep Blue supercomputer. Though Kasparov won that match, the next year Deep Blue would be the victor, marking the first time a computer program defeated a world champion in a tournament-style match. What I didn't recognize at the time was that I wasn't witnessing the beginning of a new era of AI capability, but rather the apex achievement of a rule-based era in the field that was rapidly declining. By the

time I completed research for my dissertation, both I and much of the world had moved on to a new approach to AI known as machine learning.<sup>1</sup> These transition points of different eras in the field are important for understanding both where it has been and the significance of where it may be going.

### Artificial Intelligence History

The earliest applications of AI that I demonstrated practical capabilities beyond proof-of-concepts were rule-based systems. Initially, such systems used simple rules and logic to simulate human decision-making. By the 1970s, expert systems were developed that used large sets of rules to make decisions in specialized domains. Perhaps the peak expert system was Deep Blue. With a complex but well-defined problem and sufficient computing power to search and apply rules it could outperform any human in that task. However, in general these systems were limited in their ability to handle ambiguity, learn from data, and adapt to complex problems. Their dependence on hand-crafted rules made these systems inflexible and difficult to scale.

Years later, after considerable disappointment in rule-based systems, and disinterest in AI generally, another stage of AI emerged in the 1990s with machine learning. Machine learning allowed systems to learn from data, so that

rules could be created and managed flexibly according to the data that were already collected and expected to be available in their application. To a degree, machine learning is a form of automated analysis of data, where algorithms that mimic different approaches to analyzing data are applied automatically. This made them much more flexible and useful, though they were dependent on the availability of training data and processing capabilities of the computers generating and applying the developed models. If sufficient data were not available to indicate patterns for learning, machine learning algorithms were limited in what could be discovered.<sup>2</sup>

The next major breakthrough came in the 2010s with the rise of deep learning, which is a subset of machine learning.<sup>3</sup> This was in large part due to the rapid growth of available data that could be used for machine learning and increased processing capabilities; deep learning coincided with the era of "Big Data" and the emergence of graphics processing units (GPUs) for processing of deep learning used a particular type of machine learning algorithm, artificial neural networks that were composed of multiple layers. These networks are capable of learning high-level features from raw data, which allowed greater abstraction and more powerful modeling of complex datasets.

In addition, neural network models once developed are highly efficient, because they use direct numeric calculations with the variables among the nodes in the networks. With more data and processing capabilities, deep learning allowed discovery of patterns that could not be identified during the previous machine learning era without substantial use of domain knowledge in the process.<sup>4</sup> Not surprisingly, the areas of greatest impact with deep learning were in the domains with the largest amounts of available data, including image and speech recognition, natural language processing, and autonomous systems like self-driving cars. These data sources were available due to growth of image capture, increase of text information ability, and improved data collection.

SCIENCE OF MEDICINE | FEATURE SERIES

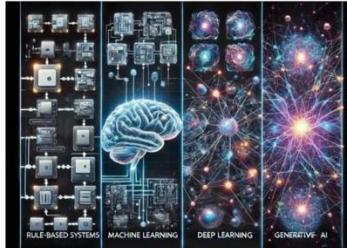


Figure 1. The different stages of artificial intelligence. This image was created (with significant prompt iteration and targeted edit requests) by DALL-E, an AI image generator.

The current era of AI, coming to prominence in the 2020s, is defined by generative AI, which extends the capabilities of AI models beyond pattern recognition and prediction to the creation of new data (Figure 1). Generative AI models such as large language models (e.g., ChatGPT) are built on transformer architectures, which are a type of deep learning model that is more efficient and powerful than previous models.<sup>5</sup> This allows high-processing capability, making generative AI algorithms the most advanced that have been applied. Importantly, generative AI provides a different kind of interaction than previous models of AI like deep learning, which were applied to specialized domains and workflows. Because generative AI creates new content rather than simply classifying content already defined, it can provide more immediate value and assist in tasks that are directly useful to a wide range of users. In fact, almost immediately after its release in late 2023, ChatGPT became the most widely and rapidly adopted technology ever.<sup>6</sup>

### Artificial Intelligence in Medicine

The application of AI to medicine has evolved alongside the broader field of AI, though progress has often been shaped by the availability of data, regulatory constraints, and the complexity of medical tasks. Expert systems were widely developed for medical settings because of their potential to standardize complex



# The role of AI in genomic medicine

- AI is likely to play an increasing role in nearly every aspect of genome-guided (and other types of) medicine
- AI has great potential to increase predictive accuracy, improve reproducibility, automate laborious tasks
- General AI is not yet attained (attainable?) - narrow AI applications will mostly augment/assist physicians and scientists in their work



Thank you!

