# Data Considerations

### **Critical Dataset Features**

Amy Abernethy, MD, PhD
Chief Medical Officer / Chief Scientific Officer & SVP - Oncology, Flatiron Health (a member of the Roche Group)
Adjunct Professor of Medicine, Duke University School of Medicine
@dramyabernethy \*\* amy@flatiron.com



### Disclosures

Amy Abernethy is Chief Medical Officer, Chief Scientific Officer and SVP, Oncology at Flatiron Health, a member of the Roche Group, and has stock ownership in Roche

### Other disclosures include:

- Athenahealth (NASDAQ: ATHN), Board of Directors, Stock ownership
- CareDx (NASDAQ: CDNA), Board of Directors, Stock ownership
- Orange Leaf Associates, LLC, Owner
- Highlander Partners, Senior Advisor (http://highlander-partners.com/)
- SignalPath Research, Advisor (http://signalpath.com/)
- The One Health Company, Special Advisor (<a href="http://www.theonehealthcompany.com/">http://www.theonehealthcompany.com/</a>)
- RobinCare, Advisor (<a href="https://getrobincare.com/">https://getrobincare.com/</a>)
- KelaHealth, Inc. peri-surgical risk prediction, Advisor (<u>kelahealth.com</u>)
- Patent Pending: Current patent pending for a technology that facilitates the extraction of unstructured information from medical records.

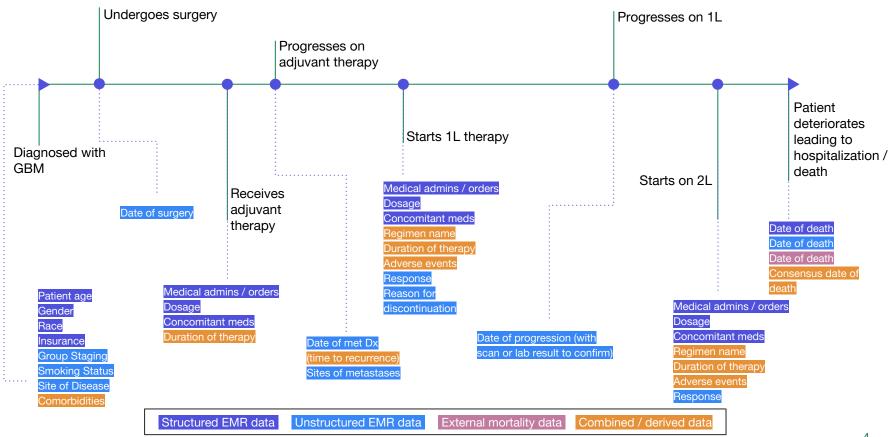


1. Importance of having all of the aspects of the data you need to create algorithms, etc

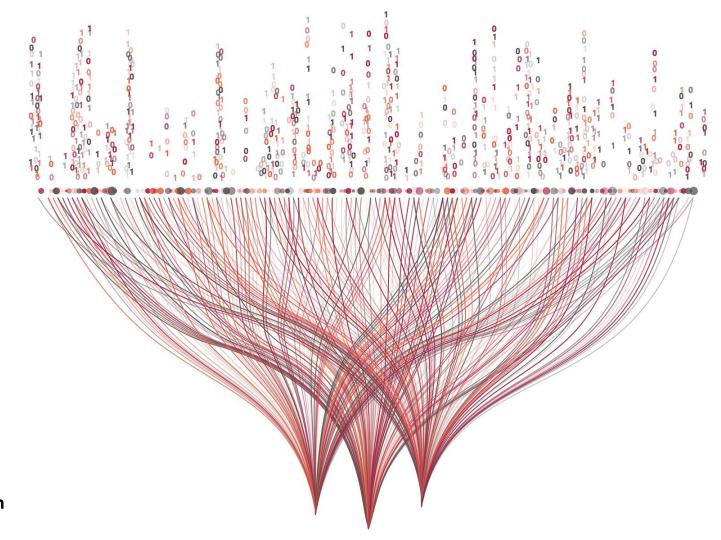
Reliable data produce reliable algorithms and reliable results



### A comprehensive view of the patient journey



\*Relative timing not exact





## Evaluating relevancy and quality

### **Data Relevancy**

- Availability of key data elements
  - Exposure
  - Outcome
  - Covariate
  - Patient-level linking (if applicable)
- Representativeness
- Sufficient subjects
- Longitudinality



### **Data Quality**

- Accuracy
  - Validity
  - Conformance
  - Plausibility
  - Consistency
- Completeness
- Provenance
- Transparency of data processing



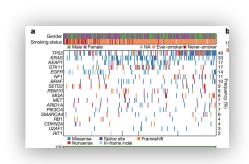
### **Fit-for-Purpose Data**

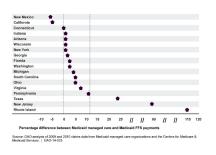
Within the given clinical and regulatory context, the real-world dataset is of sufficient quality, as well as relevant, robust, and representative.

2. Different data types have different technical, contextual and governance features which influence reliability



# Different data types have different technical, contextual and governance features









Instrumentation data (e.g., genomics, imaging, immune profiling) Administrative data (e.g., claims, mortality)

Longitudinal clinical data (e.g, electronic health records, registries)

Patient generated data (e.g., PROs, biosensors)



- Raw images are best if available, but require massive storage capabilities and are hard to de-identify
  - Pixalated data is not the same as an interpretation of findings
  - Information needed for subsequent analyses and algorithms is the interpretation
- An alternative is to focus on radiology reports in the EHR
  - Reports must be curated into a analytical dataset
  - Quality of the interpretation varies depending on reviewer
  - Findings in routine clinical practice vary from third party independent review for clincial trials
- Adequate interpretation of scans requires information about context
- Standardized systems for proxy data from scans (e.g., RECIST)
  - Rely on a series of features that may not be available outside of clinical trials (e.g., target lesions, prior scans)
  - May miss clinical features must understand the pitfalls of the approach





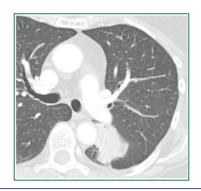
- Raw images are best if available, but require massive storage capabilities and are hard to de-identify
  - Pixalated data is not the same as an interpretation of findings
  - o Information needed for subsequent analyses and algorithms is the interpretation
- An alternative is to focus on radiology reports in the EHR
  - Reports must be curated into a analytical dataset
  - Quality of the interpretation varies depending on reviewer
  - Findings in routine clinical practice vary from third party independent review for clincial trials
- Adequate interpretation of scans requires information about context
- Standardized systems for proxy data from scans (e.g., RECIST)
  - Rely on a series of features that may not be available outside of clinical trials (e.g., target lesions, prior scans)
    - May miss clinical features must understand the pitfalls of the approach



- Raw images are best if available, but require massive storage capabilities and are hard to de-identify
  - Pixalated data is not the same as an interpretation of findings
  - o Information needed for subsequent analyses and algorithms is the interpretation
- An alternative is to focus on radiology reports in the EHR
  - Reports must be curated into a analytical dataset
  - Quality of the interpretation varies depending on reviewer
  - Findings in routine clinical practice vary from third party independent review for clincial trials
- Adequate interpretation of scans requires information about context
- Standardized systems for proxy data from scans (e.g., RECIST)
  - Rely on a series of features that may not be available outside of clinical trials (e.g., target lesions, prior scans)
    - May miss clinical features must understand the pitfalls of the approach



### Curation of Unstructured Data



#### Gross Description

The specimen is received in formalin labeled with the patient's name. It consists of a 1 x 0.3 x 0.1 cm aggregate of pink-tan to red-pink soft tissue cores and fragments entirely submitted in one block.

Dictated by: GREGORY W SMITH PA Entered: 02/06/12 - 1544 JAM

#### Microscopic Description

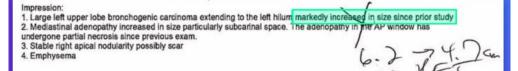
The specimen consists of a well differentiated adenocarcinoma, favor lung primary. CK7 and TTF are positive. CK20 is negative. A colleague agrees with this malignant diagnosis.

Dictated by: THOMAS J GRIFONE MD Entered: 02/07/12 - 1423 SML

#### Diagnosis

SPECIMEN SUBMITTED AS TRUCUT BIOPSY LEFT LUNG NODULE:

- WELL-DIFFERENTIATED ADENOCARCINOMA, FAVOR LUNG PRIMARY.
- SEE ABOVE.

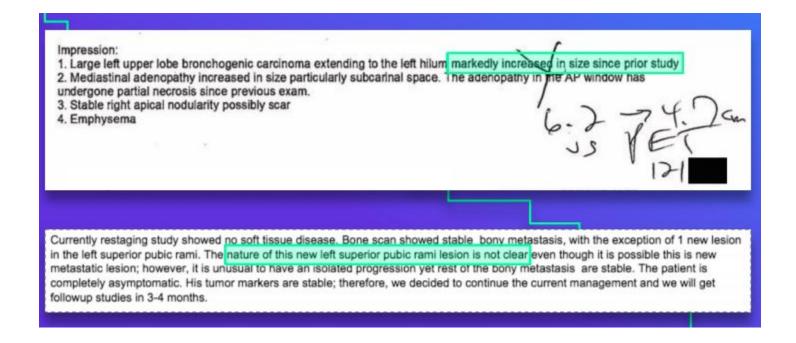


Currently restaging study showed no soft tissue disease. Bone scan showed stable bony metastasis, with the exception of 1 new lesion in the left superior public rami. The nature of this new left superior public rami lesion is not clear even though it is possible this is new metastatic lesion; however, it is unusual to have an isolated progression yet rest of the bony metastasis are stable. The patient is completely asymptomatic. His tumor markers are stable; therefore, we decided to continue the current management and we will get followup studies in 3-4 months.

- Raw images are best if available, but require massive storage capabilities and are hard to de-identify
  - Pixalated data is not the same as an interpretation of findings
  - o Information needed for subsequent analyses and algorithms is the interpretation
- An alternative is to focus on radiology reports in the EHR
  - Reports must be curated into a analytical dataset
  - Quality of the interpretation varies depending on reviewer
  - Findings in routine clinical practice vary from third party independent review for clincial trials
- Adequate interpretation of scans requires information about context
- Standardized systems for proxy data from scans (e.g., RECIST)
  - Rely on a series of features that may not be available outside of clinical trials (e.g., target lesions, prior scans)
    - May miss clinical features must understand the pitfalls of the approach



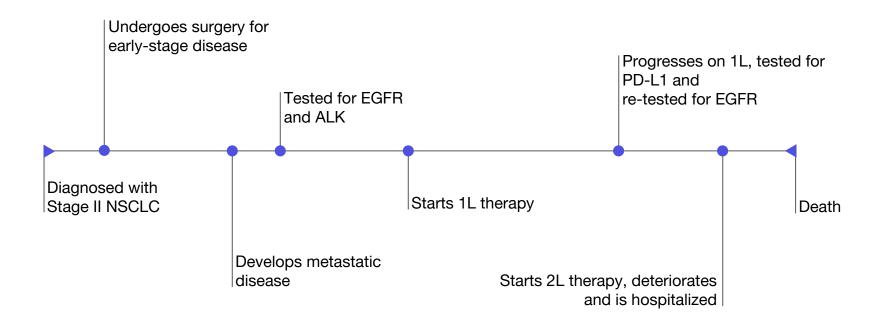
### Context: Accounting for Changing Interpretations Over Time



- Raw images are best if available, but require massive storage capabilities and are hard to de-identify
  - Pixalated data is not the same as an interpretation of findings
  - o Information needed for subsequent analyses and algorithms is the interpretation
- An alternative is to focus on radiology reports in the EHR
  - Reports must be curated into a analytical dataset
  - Quality of the interpretation varies depending on reviewer
  - Findings in routine clinical practice vary from third party independent review for clincial trials
- Adequate interpretation of scans requires information about context
- Standardized systems for proxy data from scans (e.g., RECIST)
  - Rely on a series of features that may not be available outside of clinical trials (e.g., target lesions, prior scans)
    - May miss clinical features must understand the pitfalls of the approach







Clinical events are a combination of clinical, pathological, radiological, & biomarker data - *in context* 



Undergoes surgery for early-stage disease Progresses on 1L, tested for PD-L1 and Tested for EGFR re-tested for EGFR and ALK **Diagnosed with** Starts 1L therapy Stage II NSCLC Death Develops metastatic disease Starts 2L therapy, deteriorates and is hospitalized Gross Description The specimen is received in formalin labeled with the patient's name. It consists of a 1 x 0.3 x 0.1 cm aggregate of pink-tan to red-pink soft tissue cores and fragments entirely submitted in one block. Dictated by: GREGORY W SMITH PA Entered: 02/06/12 - 1544 JAM Microscopic Description The specimen consists of a well differentiated adenocarcinoma, favor lung primary. CK7 and TTF are positive. CK20 is negative. A colleague agrees with this malignant diagnosis. Dictated by: THOMAS J GRIFONE MD Entered: 02/07/12 - 1423 SML SPECIMEN SUBMITTED AS TRUCUT BIOPSY LEFT LUNG NODULE:

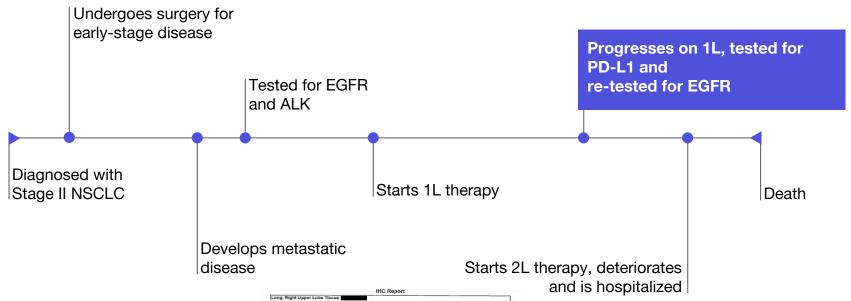
- NELL-DIFFERENTIATED ADENOCARCINOMA, FAVOR LUNG PRIMARY.

- SEE ABOVE.

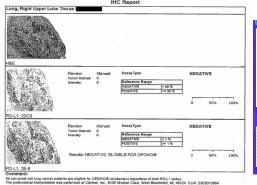


Undergoes surgery for early-stage disease Progresses on 1L, tested for PD-L1 and Tested for EGFR re-tested for EGFR and ALK Diagnosed with Starts 1L therapy Stage II NSCLC Death **Develops metastatic** disease Starts 2L therapy, deteriorates and is hospitalized Path?





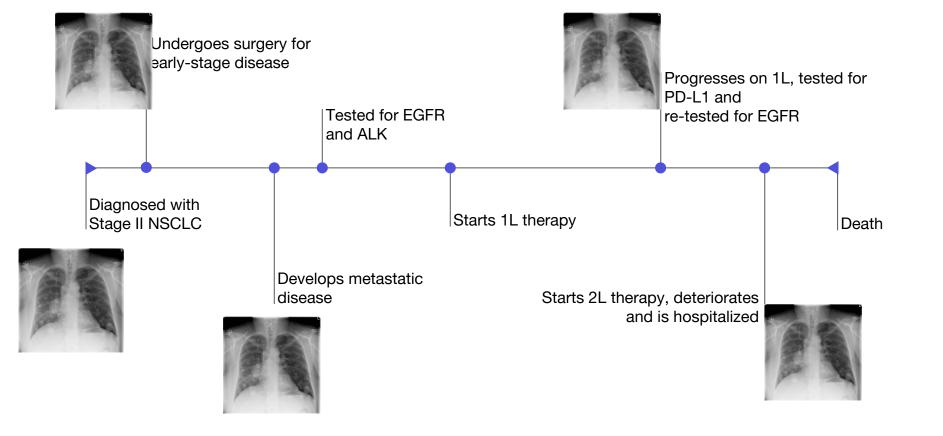
Time to progression is dependent on when patient is evaluated





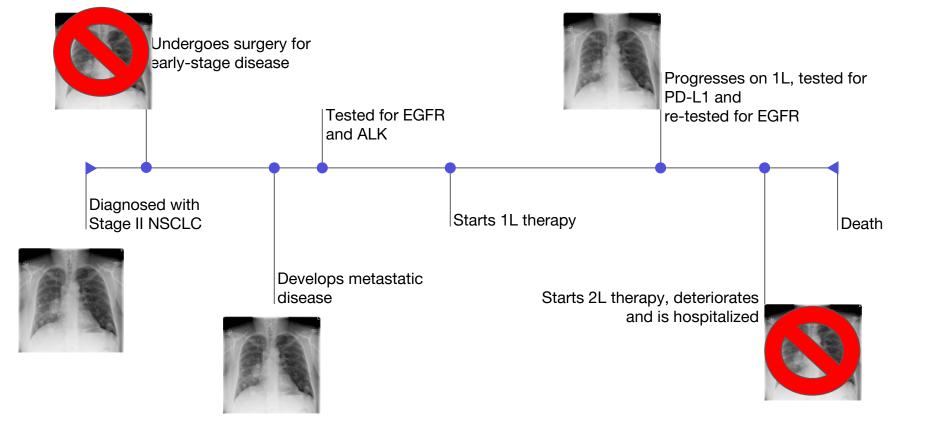
Currently restaging study showed no soft lissue disease. Bone scan showed stable, bony metastasis, with the exception of 1 new lesion in the left superior pubic ram: lesion is not clear/even though it is possible this is new metastatic lesion; however, it is unusual to have an solised progression yet rest of the cony metastasis are stable. The patient is completely asymptomatic. His tumor markers are stable; therefore, we decided to continue the current management and we will get followup studies in 3-4 months.





## Complete datasets increase reliability





## Complete datasets increase reliability



- Raw images are best if available, but require massive storage capabilities and are hard to de-identify
  - Pixalated data is not the same as an interpretation of findings
  - o Information needed for subsequent analyses and algorithms is the interpretation
- An alternative is to focus on radiology reports in the EHR
  - Reports must be curated into a analytical dataset
  - Quality of the interpretation varies depending on reviewer
  - Findings in routine clinical practice vary from third party independent review for clincial trials
- Adequate interpretation of scans requires information about context
- Standardized systems for proxy data from scans (e.g., RECIST)
  - Rely on a series of features that may not be available outside of clinical trials (e.g., target lesions, prior scans)
    - May miss clinical features must understand the pitfalls of the approach



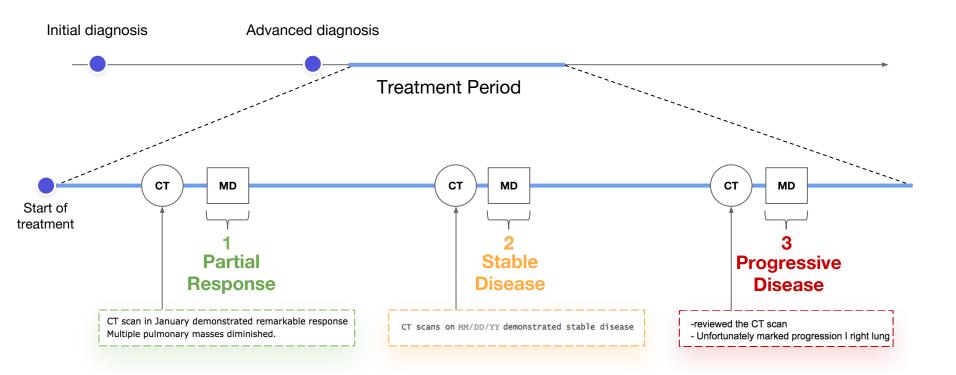


# rwTR Assessment Categories

Category	Description
Complete Response (CR)	Complete resolution of all visible disease.
Partial Response (PR)	Partial reduction in size of visible disease in some or all areas without any areas of increase in visible disease. Captures a decrease in disease volume even though disease is still present.
Stable Disease (SD)	No change in overall size of visible disease. Also includes cases where some lesions increased in size and some lesions decreased in size.
Progressive Disease (PD)	Increase in visible disease and/or presence of any new lesions. Includes cases where the clinician indicates progressive disease, PD, or POD as the overall assessment.
Pseudoprogression	Clinician indicates pseudoprogression or related terminology (e.g., tumor flare).
Indeterminate response	Clinician specifically indicates that the response is "indeterminate" or "uncertain," or if the clinician's interpretation of the scan(s) cannot be mapped to one of the above assessment categories.
Not documented	Clinician's note references this response imaging (e.g., "Patient had recent scan") but does not mention any assessment of tumor response.



# Response is tied to exposure to a therapy





### Governance and De-ID



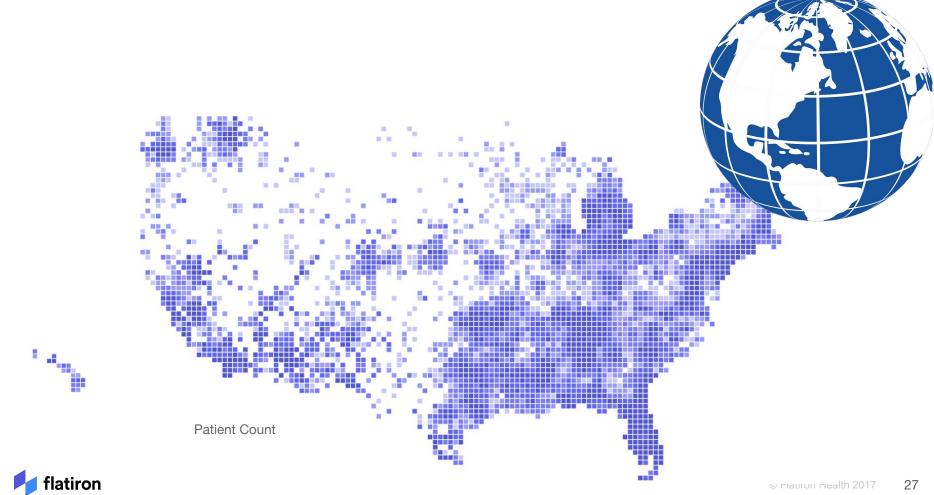
- Who approves access?
- What kind of access?
- How is access QA'ed?

All ultrasound machines store the patient's name and MRN as intrinsic parts of the image and are displayed whenever the image is displayed (<u>Fig. 3</u>). The DICOM headers can be deidentified, but the PHI in the image remains. In addition to the patient's name and MRN, some formats display the birth date, scan date and time, or facility name. https://www.ajronline.org/doi/full/10.2214/AJR.13.11789

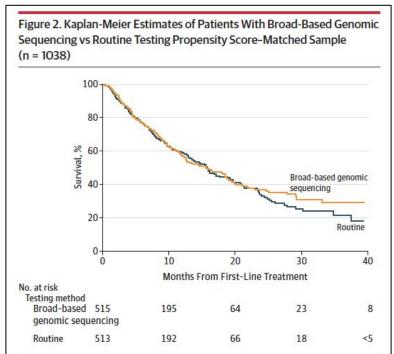


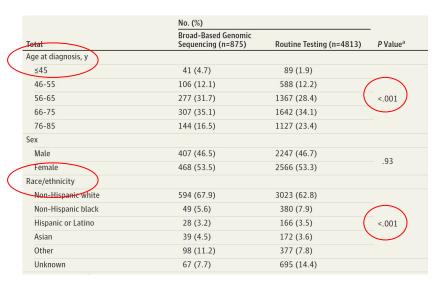
3. Representativeness of datasets (genomic, geographic); risk of bias in the underlying data





# Does genomic testing improve survival for lung cancer patients?





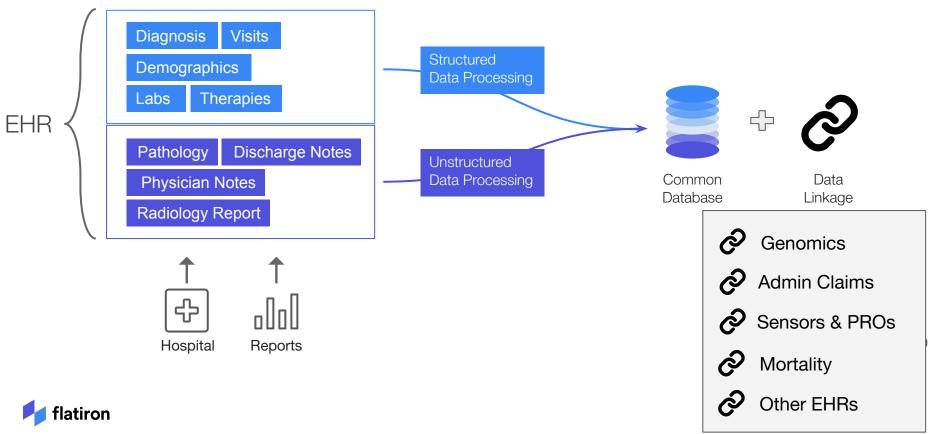
Presley et al. JAMA 2018



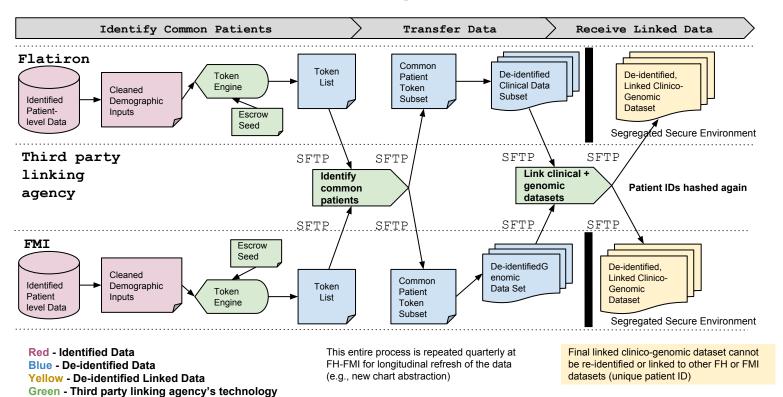
# 4. Methods to resolve data gaps



# Data linkage to fill gaps



# Data linkage solutions

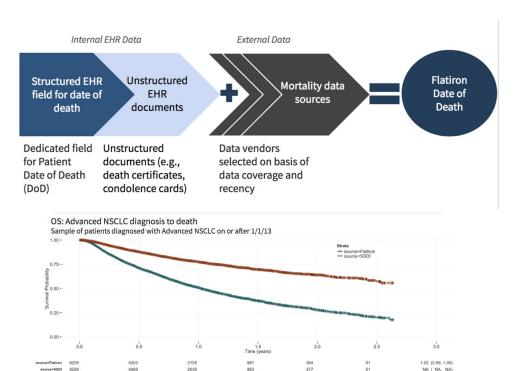


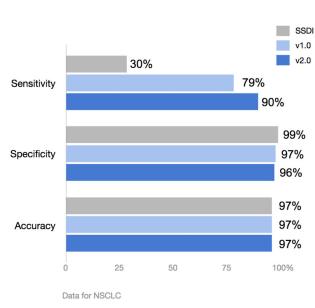




# Evaluate data against a reference standard

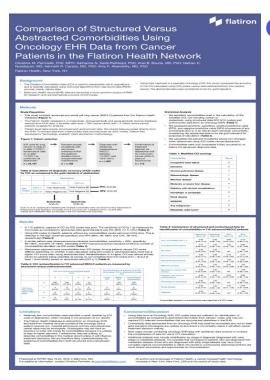
E.g., gold standard = National Death Index







## Developing proxy and surrogate variables when needed



# Comparison of structured vs. abstracted comorbidities Results:

- Capture of CCI by ICD codes was poor. The sensitivity of CCI ≥ 1 as measured by ICD codes as compared to abstracted data (gold standard) was 9% (95% CI: 5-14%)
- Using ICD codes to identify patients without any comorbidities works well most of the time. This is reflected in the high observed specificity and NPV (99%, CI: 98-100% and 72%, CI: 69-75%, respectively)
- Abstraction captured more comorbidities than ICD codes. However, comorbidity data abstracted from an Oncology EHR may itself be incomplete and not an ideal gold standard. Oncologists are unlikely to document a comorbidity unless it will affect cancer treatment decision-making

5. Issues in Precision Medicine: Available biologic/biomarker data, small cohorts



# Deep Phenotype is Needed

In an era of abundant data, merging biological information with deep clinical phenotype is more important than ever



The Human Phenotype Ontology project: linking molecular biology and disease through phenotype data

phenotype data

# The details of disease

Precision medicine demands precise matching of deep genomic and phenotypic models — and the deeper you go, the more you know.

S14 | NATURE | VOL 527 | 5 NOVEMBER 2015





Nucleic Acids Research, 2013, 1-9 doi:10.1093/nar/gkt1026

# Precision medicine often yields small cohorts

The NEW ENGLAND JOURNAL of MEDICINE

#### ORIGINAL ARTICLE

### Vemurafenib in Multiple Nonmelanoma Cancers with BRAF V600 Mutations

David M. Hyman, M.D., Igor Puzanov, M.D., Vivek Subbiah, M.D.,
Jason E. Faris, M.D., Ian Chau, M.D., Jean-Yves Blay, M.D., Ph.D.,
Jürgen Wolf, M.D., Ph.D., Noopur S. Raje, M.D., Eli L. Diamond, M.D.,
Antoine Hollebecque, M.D., Radj Gervais, M.D.,
Maria Elena Elez-Fernandez, M.D., Antoine Italiano, M.D., Ph.D.,
Ralf-Dieter Hofheinz, M.D., Manuel Hidalgo, M.D., Ph.D.,
Emily Chan, M.D., Ph.D., Martin Schuler, M.D., Susan Frances Lasserre, M.Sc.,
Martina Makrutzki, M.D., Florin Sirzen, M.D., Ph.D., Maria Luisa Veronese, M.D.,
Josep Tabernero, M.D., Ph.D., and José Baselga, M.D., Ph.D.

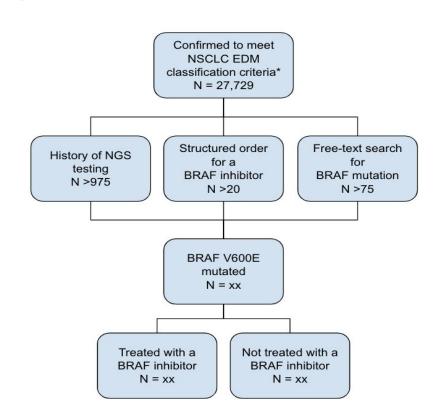
#### ABSTRACT

#### BACKGROUND

BRAF V600 mutations occur in various nonmelanoma cancers. We undertook a histology-independent phase 2 "basket" study of vemurafenib in BRAF V600 mutation—positive nonmelanoma cancers.

#### METHODS

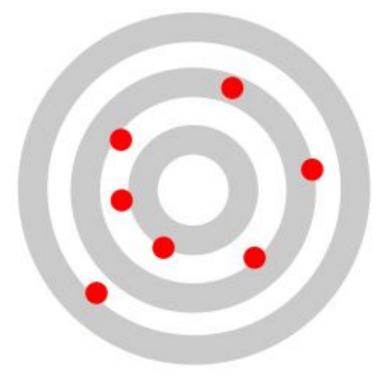
We enrolled nationts in six prespecified capper soborts, nationts with all other



# 6. Data reliability and quality

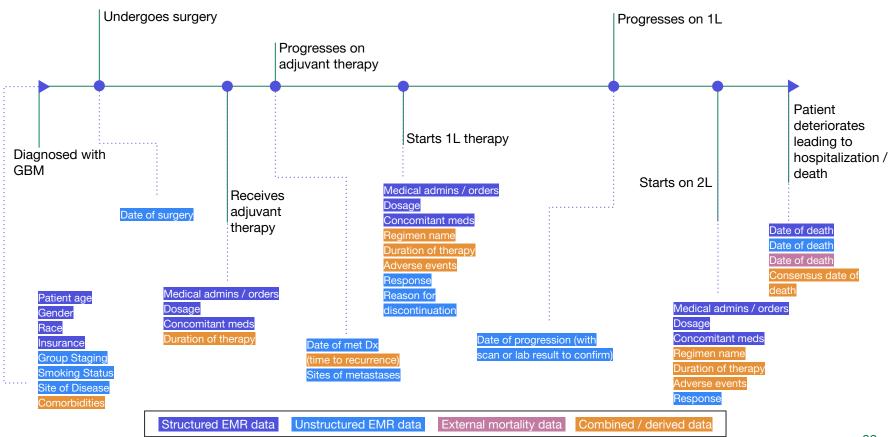


# Reliable data produce reliable algorithms and reliable results





### Documentation of source, quality and provenance



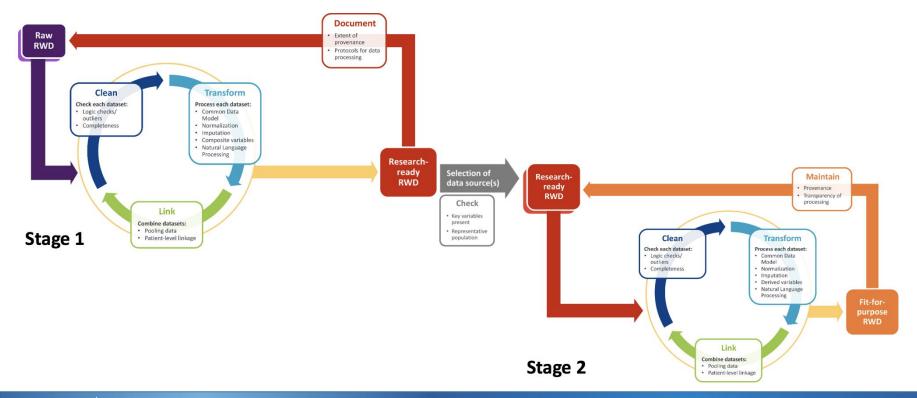
\*Relative timing not exact

# Need a consistent approach to documenting completeness and quality

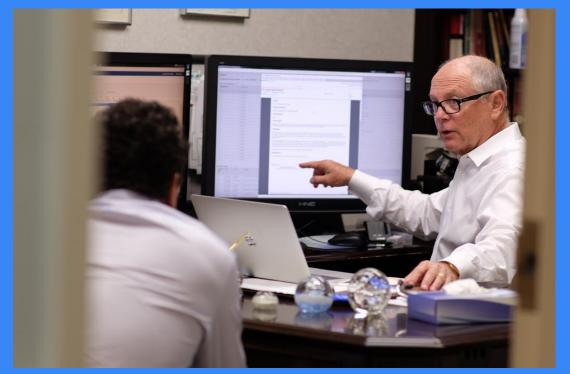
Project:	FDA						
						Kappas scale	
Note: For questions where a high percentage of patients have a common answer (e.g., PD-L1 testing status), kappa may be significantly lower than inter-rater agreement.						0.8 to 1.0	
In these cases, it may be more accurate to use inter-rater agreement to measure reliability of the data.						0.6 to 0.8	
					Moderate	0.4 to 0.6	
					Fair	0.2 to 0.4	
					Slight	0 to 0.2	
Table: Enhanced_Advance	dNSCLC						
Summary of variable inter-	rater agreement and kappas						
		Corresponding question(s) on abstraction		Inter-rater agreement	Kanna	Kappa (30-day	
Variable	Description of variable	form	Question type	(exact day for dates)	(exact agreement)	window for dates	
DiagnosisDate	Date of initial diagnosis	Enter the date of initial diagnosis	date	0.795	0.794	0.902	
		Enter the date of the first		0.730	0.734	0.302	
AdvancedDiagnosisDate	Date of diagnosis of advanced disease: first recurrence or metastasis	diagnosis of metastatic or advanced NSCLC	date	0.695	0.695	0.796	
MetastaticDiagnosisDate		Enter the date of initial					
		diagnosis [for ~55% of patients in the cohort who	200000				
			date	0.795	0.794	0.902	
WetastaticDiagnosisDate	Date of diagnosis of metastatic disease	Enter the date of distant metastatic diagnosis [for					
		~45% of patients in the cohort who are diagnosed					
		non-metastatic]	date	0.527	0.476	0.557	
Histology	Histology	Select the histology	drop down	0.947	0.894		
GroupStage	Group stage at time of initial diagnosis	Select the group stage	drop down	0.848	0.768		
SmokingStatus	Documented history of smoking	Smoking status Was the tumor tested for a	drop down	0.934	0.695		
EgfrTested	Indicator of whether the tumor was tested for a EGFR mutation	EGFR mutation?	boolean	0.927	0.84		
AlkTested	Indicator of whether the tumor was tested for an ALK rearrangement	Was the tumor tested for an ALK rearrangement?	boolean	0.901	0.791		
PdL1Tested	Indicator of whether the tumor was tested for PD-L1 expression	Was the tumor tested for PD-L1 expression?	boolean				
rutiresleu	indicator of whether the turnor was tested for FD-L1 expression	Was the tumor tested for a	DOOLGAII	0.901	0.547		
KrasTested	Indicator of whether the tumor was tested for a KRAS mutation	KRAS mutation?	boolean	0.894	0.728		
	Indicator of whether the tumor was tested for a ROS-1	Was the tumor tested for a		1		1	



# Research-ready databases



# 7. Recency and availability





# 8. Privacy and security





# 9. Policy Considerations

# Reliable data produce reliable algorithms and reliable results



# Policy Considerations

- Require accurate complete datasets that incorporate many data types
- Details about the data matter
- Biased and unreliable data lead to biased and unreliable algorithms
- Accessible data at point of care
- While maintaining patient privacy
- Solving this requires input of many types of experts - clinical, analytic, software, hardware, privacy, etc

- Stimulate tech innovation that solves all of these problems - do not skimp
- Expect transparency about data reliability and quality - as well as accuracy of output from algorithms
- Create standards for documenting reliability, quality and accuracy at the data source, dataset and algorithm level
- Develop a the workforce

