

Digital Transformation of Healthcare

ETL & Assessing Data Quality

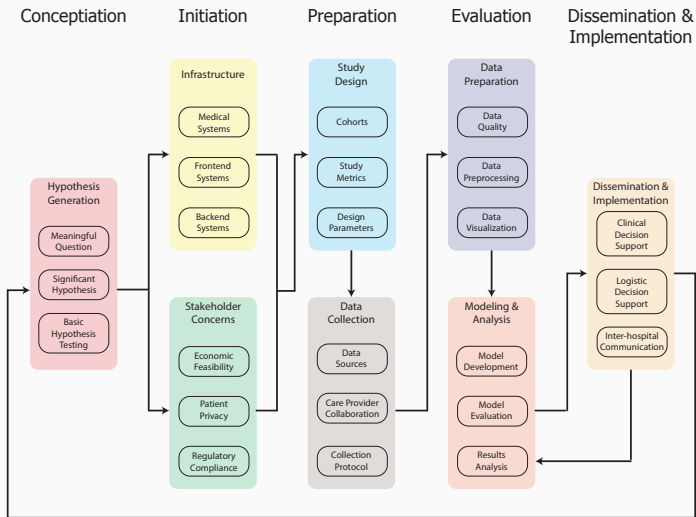
Michael Snow, M.D. Ph.D., Glen Ferguson, Ph.D.

Center for Health Data Innovations

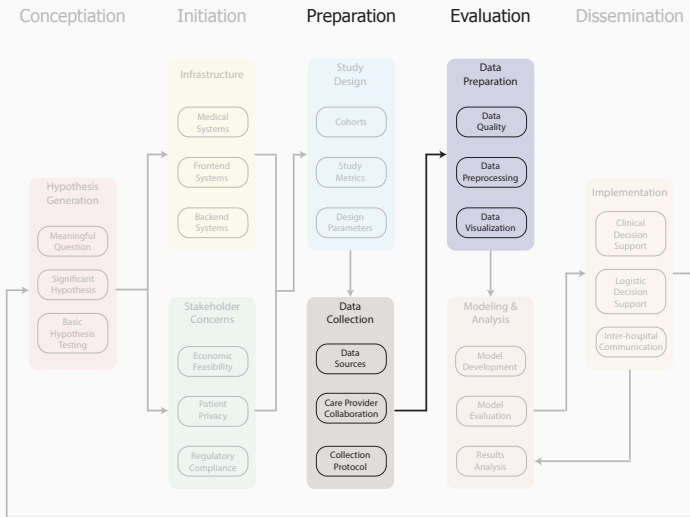
After this lecture students will be able to

- Describe the components of an ETL pipeline
- Extract variables from data and incorporate additional information
- Assess the quality of data
- Trace the steps where data quality can be affected
- Examine data for problems and discuss possible causes
- Design a process for the imputation of missing data

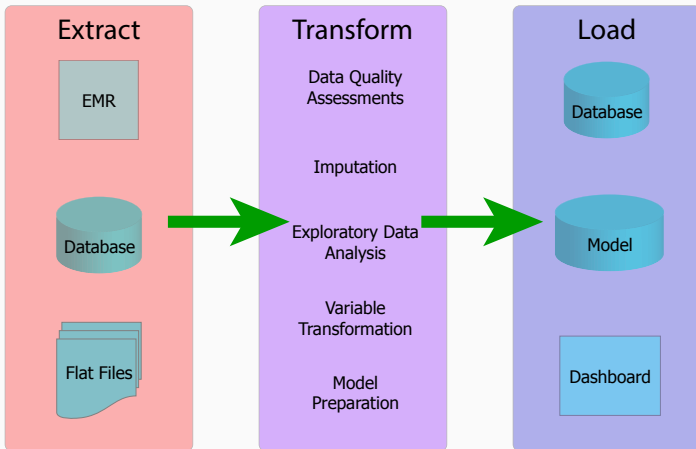
Bioinformatics Pipeline



ETL & Data Quality



Extract, Transform and Load (ETL)



Analysis is only ever as good as the data it's built upon.

- What is data quality? What makes data high quality vs low quality?
- Where along the process can you affect data quality?
- How can you design a study to collect high quality data (Quality assurance)?
- How can you identify and correct errors during and after data collection (Quality control)?

- DICOM - Digital Imaging and Communications in Medicine - is the international standard for medical images and related information. It defines the formats for medical images that can be exchanged with the data and quality necessary for clinical use
- DICOM groups information into data sets, e.g., an x-ray would contain the patient ID within the file, so that the image can never be separated from this information by mistake.
- DICOM Value Representations

<https://www.dicomstandard.org/about/>

Quality Assurance - DICOM

name	VR	value
Group Length	UL	532
Image Type	CS	DERIVED
SOP Class UID	UI	1.2.840.10008.5.1.4.1.1.2
SOP Instance UID	UI	1.2.840.114356.2008.11.30.12.34.2.329.999
Study Date	DA	20081230
Content Date	DA	20081230
Study Time	TM	122731
Content Time	TM	12299.0000
Modality	CS	CT
Institution Name	LO	Manhasset Diagnostic Imaging
Station Name	SH	
Study Description	LO	MOSES CT Outside Reference Images
Procedure Code Sequence	SQ	[{(0008, 0100): (0008, 0100) Code Value ...
Code Value	SH	MOSESOUTREFCT
Coding Scheme Designator	SH	GEIIS
Coding Scheme Version	SH	0
Code Meaning	LO	MOSES CT Outside Reference Images
Series Description	LO	Reformatted
Referenced SOP Class UID	UI	1.2.840.113619.2.51762891606.1649.1005918257.250
Referenced SOP Instance UID	UI	1.2.840.114356.2008.11.30.12.34.2.329.1301

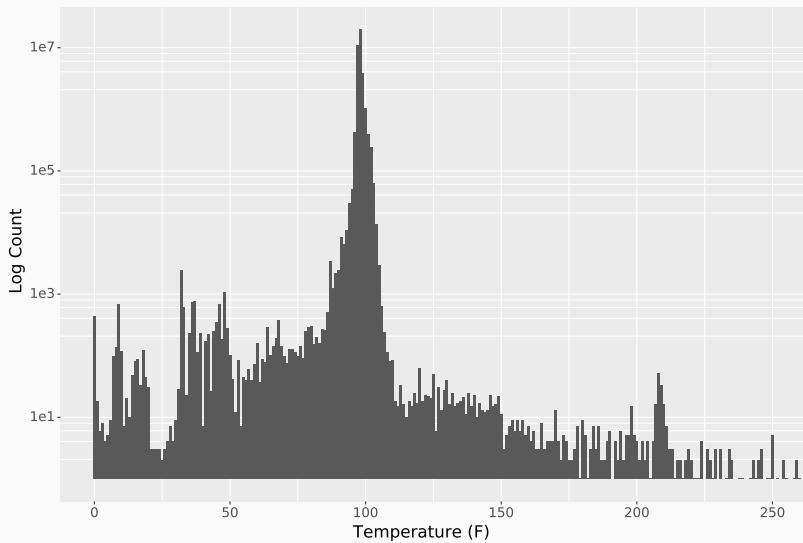
name	VR	value
Study Date	DA	20081230
Content Date	DA	20081230
Study Time	TM	122731
Content Time	TM	12299.0000

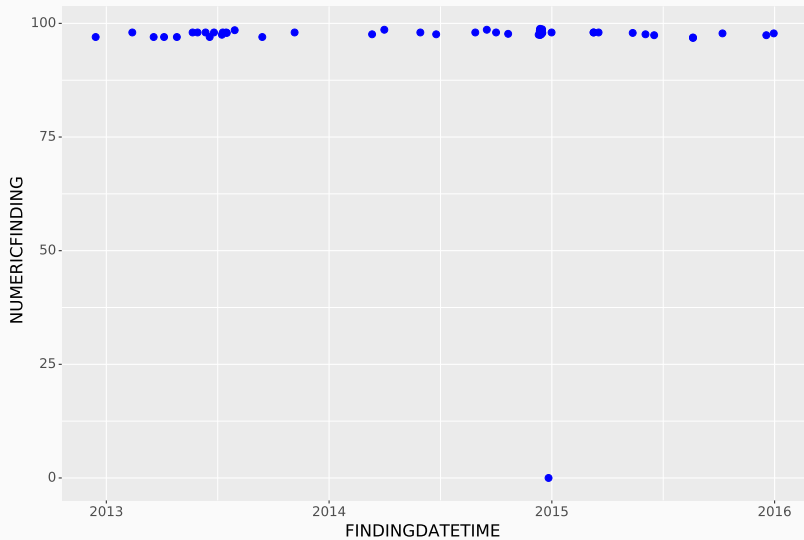
- **DA** - A string of characters of the format YYYYMMDD
- **TM** - A string of characters of the format HHMMSS.FFFFFFFF.
 - One or more of the components MM, SS, or FFFFFFFF may be unspecified as long as every component to the right of an unspecified component is also unspecified

Whose fault is this?

How can I find the temperatures recorded
from every patient in the hospital

To the SQL





Associated Values

FINDINGDATETIME	FINDINGDESC	NUMERICFINDING
2014-12-26	PULSE OXIMETRY	97.00
2014-12-26	WEIGHT/SCALE (ounces)	2800.16
2014-12-26	HEIGHT (inches)	62.00
2014-12-26	Diastolic Blood Pressure	82.00
2014-12-26	Systolic Blood Pressure	139.00
2014-12-26	HEIGHT (CM)	157.48
2014-12-26	PULSE	75.00
2014-12-26	BODY MASS INDEX	32.13
2014-12-26	O2 SAT%	97.00
2014-12-26	TEMPERATURE (F)	0.00
2014-12-26	Systolic Blood Pressure	139.00
2014-12-26	WEIGHT (KG)	79.38
2014-12-26	Diastolic Blood Pressure	82.00

Can we develop a systematic way to deal with missing data

- What are the different ways that data could be missing

- WHO data quality
- Healthcare Data Warehousing and Quality Assurance
- (2002). Defining and improving data quality in medical registries
JAMIA, 9(6), 600-611.