# How to Trust, but Verify, in Healthcare
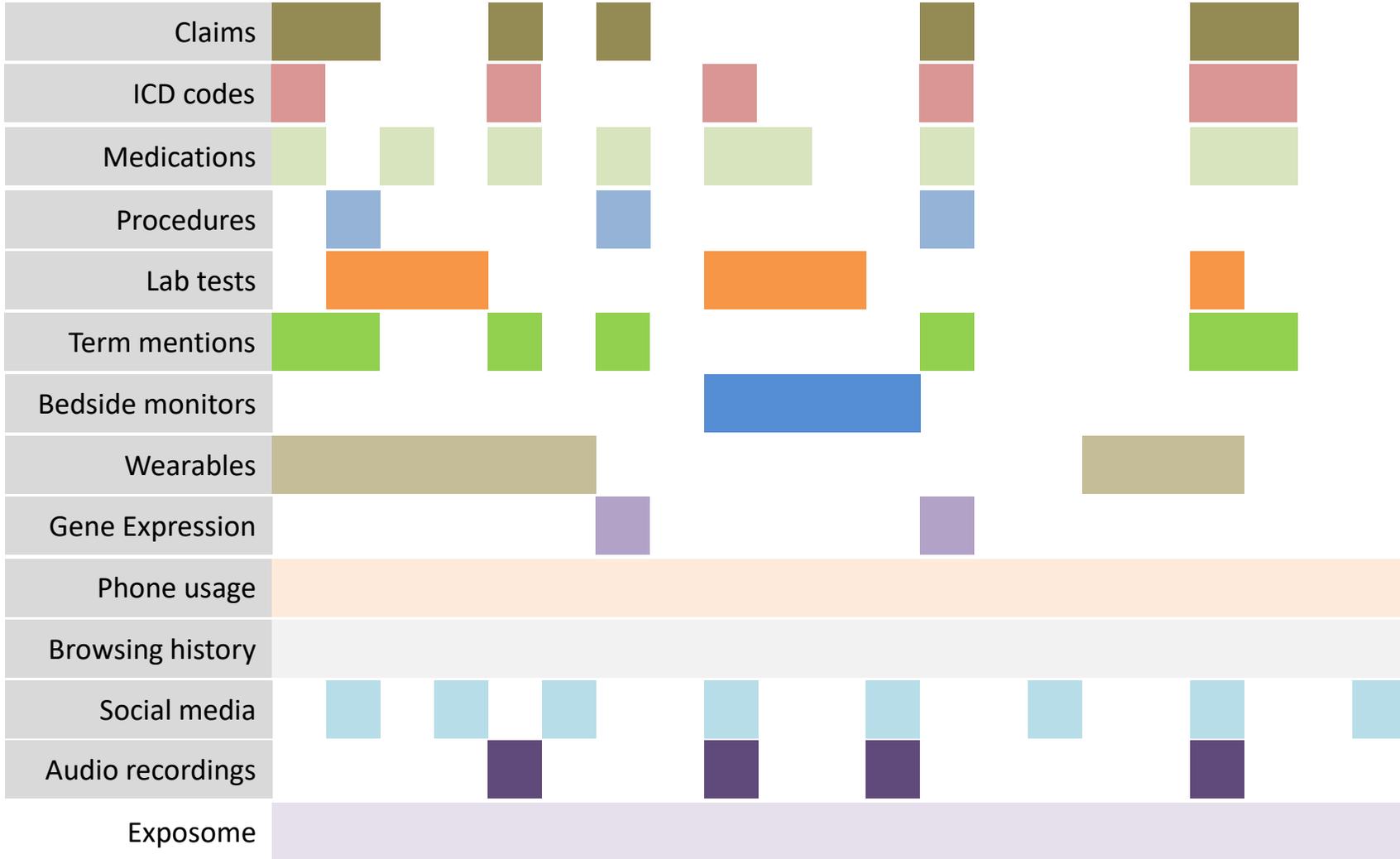
Nigam Shah, MBBS, PhD

nigam@stanford.edu

Patient Journey

Claims

ICD codes

Medications

Procedures

Lab tests

Term mentions

Bedside monitors

Wearables

Gene Expression

Phone usage

Browsing history

Social media

Audio recordings

Exposome

:

# A model

# A workflow

# Inpatient Hospital Medicine Initiated Workflow – Goals of Care Workflow

**Summary Stats**

Total Steps: 21
Level 1: 7
Level 2: 7
Level 3: 7

Handoffs: 7

Total Cycle Time: 48 hours

## CNS

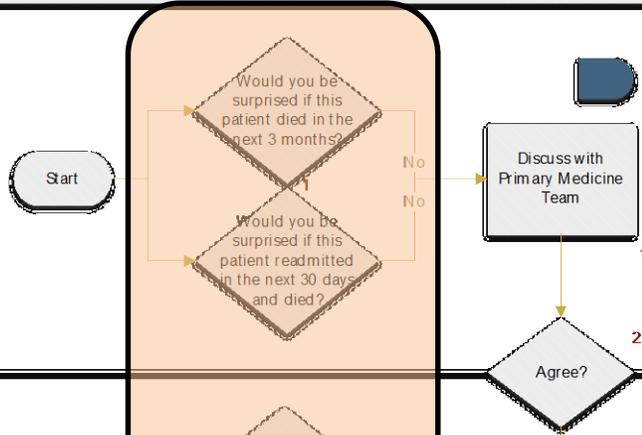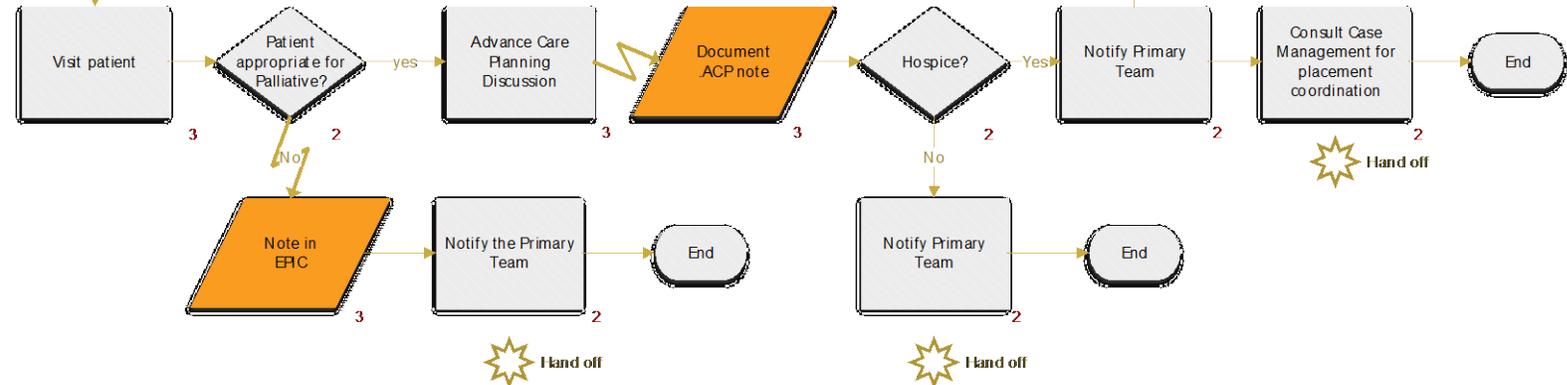Start → **Would you be surprised if this patient died in the next 3 months?** / **Would you be surprised if this patient readmitted in the next 30 days and died?** — No / No → Discuss with Primary Medicine Team [1]

Hand off

Agree? [2]

## Hospital Med Attending

Start → **Would you be surprised if this patient died in the next 3 months?** / **Would you be surprised if this patient readmitted in the next 30 days and died?** — No / No → Contact Primary Care Provider, Caregiver, and/or Family member [1] → Consult Palliative? [1] — No → Have Goals of Care discussion with patient including care giver or family [1] → Document GOC note [3] → Hospice? [1] — Yes → Place on Comfort Care Measures [3] → Consult Case Management for placement coordination [3] → End

Hospice? — No → Cont. Trmt

Consult Palliative? — Yes → Hand off

Hand off

Hand off

## Palliative Medicine

Visit patient [3] → Patient appropriate for Palliative? — yes [2] → Advance Care Planning Discussion [3] → Document ACP note [3] → Hospice? [2] — Yes → Notify Primary Team [2] → Consult Case Management for placement coordination [2] → End

Patient appropriate for Palliative? — No → Note in EPIC [3] → Notify the Primary Team [2] → End

Hospice? — No → Notify Primary Team [2] → End

Hand off

Hand off

Hand off

**Reliability Level:**
(1) Individuals: Feedback, checklists, training, basic standards
(2) Procedures: Embedded standard work, reminders, constraints
(3) Systems: fail safes, physical layout, built-in feedback, automated systems, concentration of responsibility

Margaret Smith, MBA, Value Improvement, SHC

# Definitions and Clarifications

- Trustworthiness: of the model, or the workflow around it, or both?

- Trust = proof over time that a thing does what it claims to do. Trust is earned [over time].

- HOW = interpretability
- WHY = explainability

# When predicting 24 hr. mortality …

- Interpretability is a poor surrogate for trust
  - Knowing 'how' does not help you decide what action to take

- Explainability is a poor surrogate for trust
  - Knowing 'why' does not help you decide what action to take

- Knowing that the model's prediction has helped make good decisions in the past 2 years.

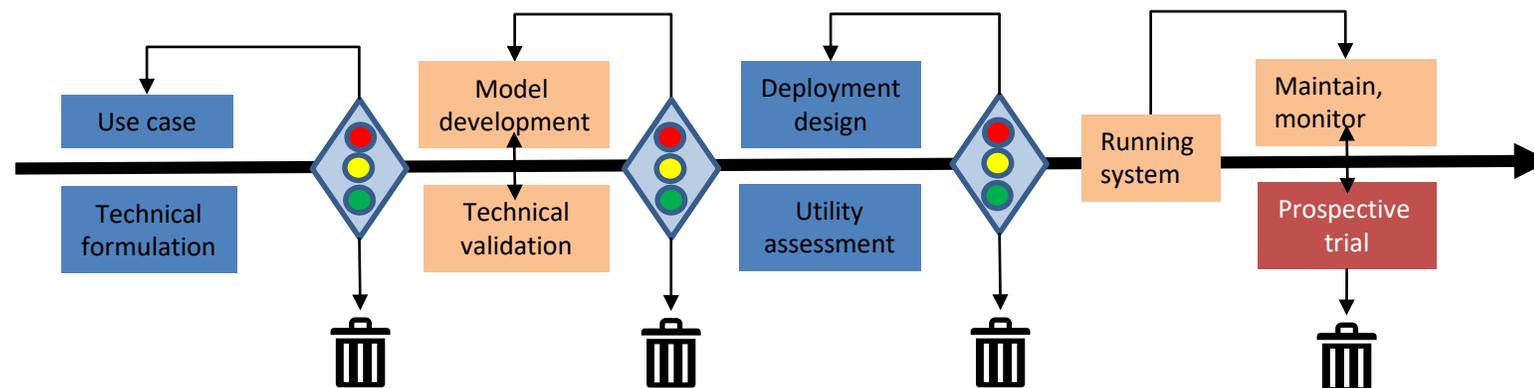# Building trustworthy (and useful!) models

**How do we get the best f: X -> Y?**
  Does representation learning help?
  Does multi-task learning help?
  Does using textual content help?
  How do we train fair models?

**Can we use f: X -> Y in the real world?**
  Can we get the data by 5 am, to make prediction by 6 am?

**Running system = model applied to each case + execution of workflow.**
  • Evaluate the impact of the *running system* on the outcomes we care about
  • Maintenance is huge liability – who will carry the pager?
  • Monitoring is unexplored



**Use case**
  What clinical outcome(s) are you trying to affect?
  Who is the target population?
  What action would you take?
  Who will take that action?

**f: X -> Y subject to...**
  use an existing equation vs. learn a new equation.

**Utility assessment**
  Given the costs of the actions and its benefit, is there net utility?

**Deployment design**
  Do we increase the efficiency of existing workflows
  Do we require entirely new workflows

# Acknowledgements