

What Can Explainability Mean for Probabilistic Inference?

Human-aware PGMs and Probabilistic Inference via Lifted Model Reconciliation – A KI Starter Project Starting in 2024

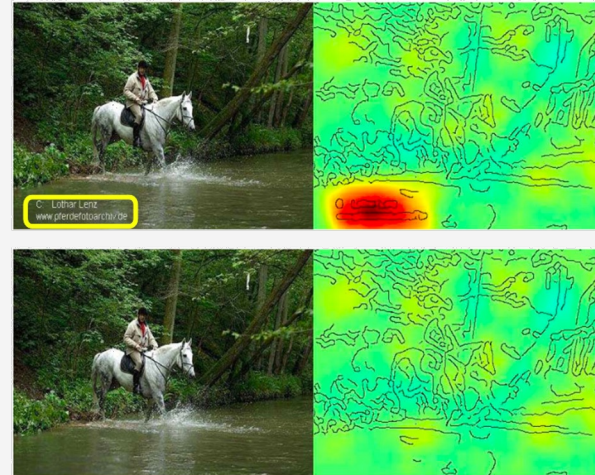
Tanya Braun



Explainable AI (XAI)

- Explanations critical for trust, collaboration, ...
- *Explain classifications*
 - Debugging tool for inscrutable representations
 - “Pointing” explanations

Horse-picture from Pascal VOC data set



Artificial picture of a car

Source tag present

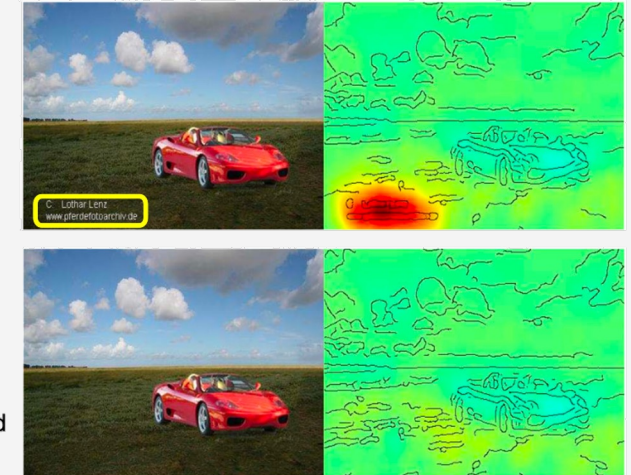


Classified as horse

No source tag present



Not classified as horse



Prediction:
School bus

Difference between left
and right magnified by 10

Prediction:
Ostrich

Please point to
the “ostrich” part

Explainable AI (XAI)

- *Explain decisions / plans*
 - Classical planning: Given a planning problem (Σ, s_0, S_g) (agent model \mathcal{M}^R)
 - Find a plan $\pi = \langle a_1, a_2, \dots, a_n \rangle$ that transforms s_0 to a state $s_n \in S_g$
 - Human-aware planning: Team of human and agent (e.g, robot)
 - Agent's model \mathcal{M}_r^H of human's model \mathcal{M}^H
 - Allows the agent to **anticipate human behaviour** to assist, avoid, team
 - Agent's model \mathcal{M}_h^R that the agent expects the human to have of \mathcal{M}^R
 - Allows the agent to **anticipate human expectations** to conform to those expectations, explain its own behaviour in terms of those expectations
 - Interpretability: Set up explicable / legible / predictable / ... plans
 - Explanations: *Reconcile* model differences

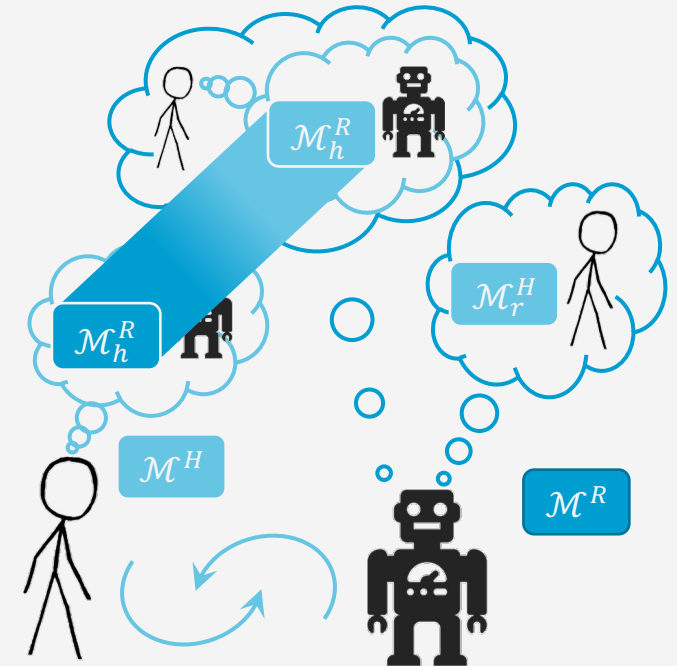
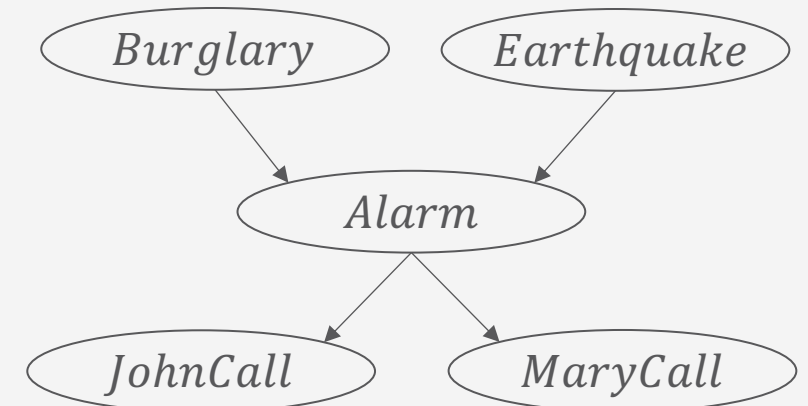
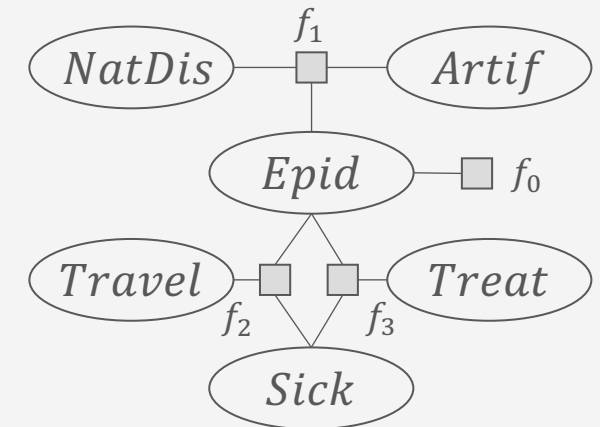


Figure based on Kambhampati

What about Probabilistic Inference?

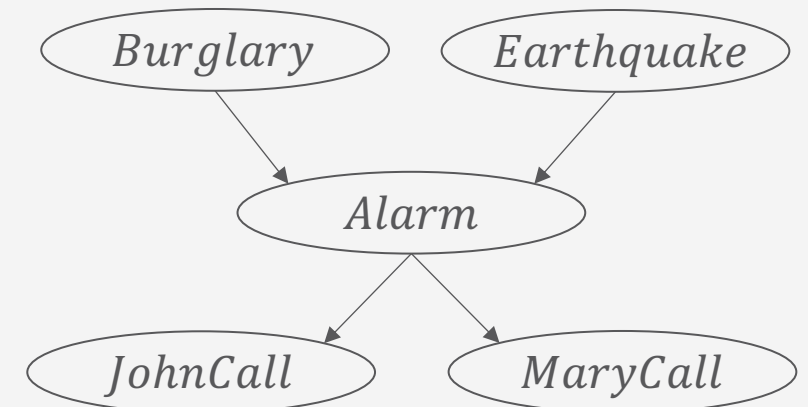
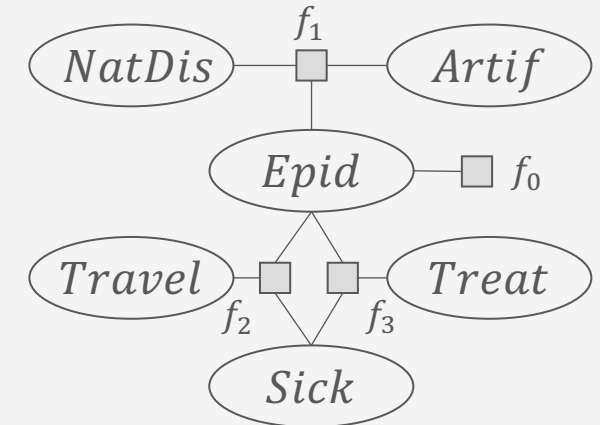
- Probabilistic graphical models (PGMs)
 - Set of random variables, set of conditional probability distributions / factors, connecting random variables
 - E.g., Bayesian network (above), factor graph (below)
 - Semantics: full joint probability distribution as (normalised) product of distributions / factors
- Probabilistic inference
 - Query for probability (distribution) of event / random variable
 - E.g.,
 - $P(\text{Burglary} = \text{true} | \text{MaryCall} = \text{true})$
 - $P(\text{Epid} | \text{Sick} = \text{true})$
 - Solve by eliminating all non-query terms (sum-out)



What Can We Find in the Direction of XAI?

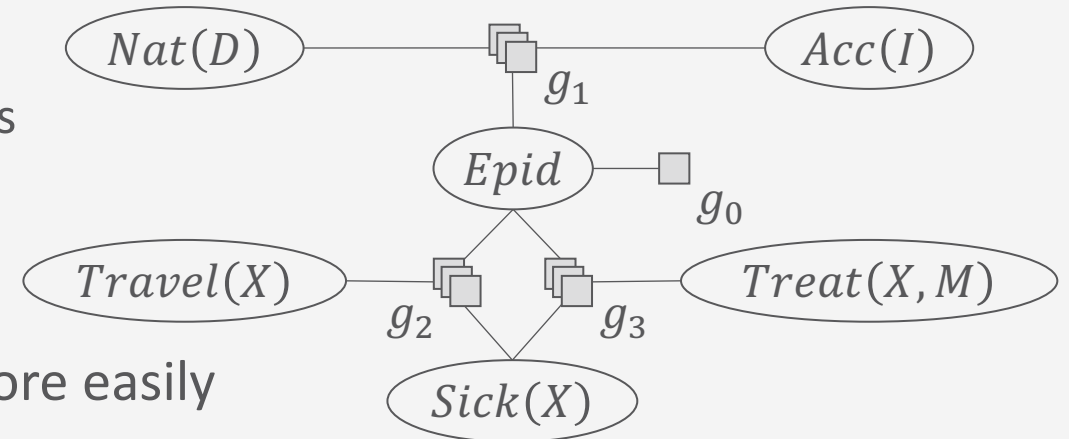
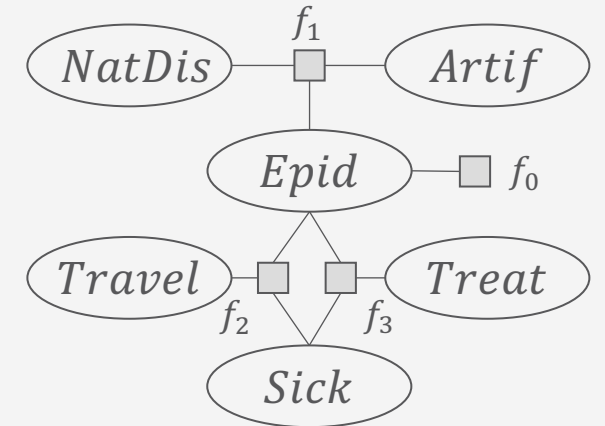
- Most Probable Explanation (MPE)
 - Given a set of events, what is the most probable state of the remaining random variables?
 - Formally, given a PGM G and a set of events e :

$$MPE_G(e) = \arg \max_{v \in \text{Val}(V)} P(v \mid e) = \arg \max_{v \in \text{Val}(V)} P(v, e)$$
 - $V = \text{rv}(G) \setminus \text{rv}(e)$ random variables of G without those in e
 - Max-out instead of sum-out to answer an MPE query instead of a probability query
- Is that really explaining?



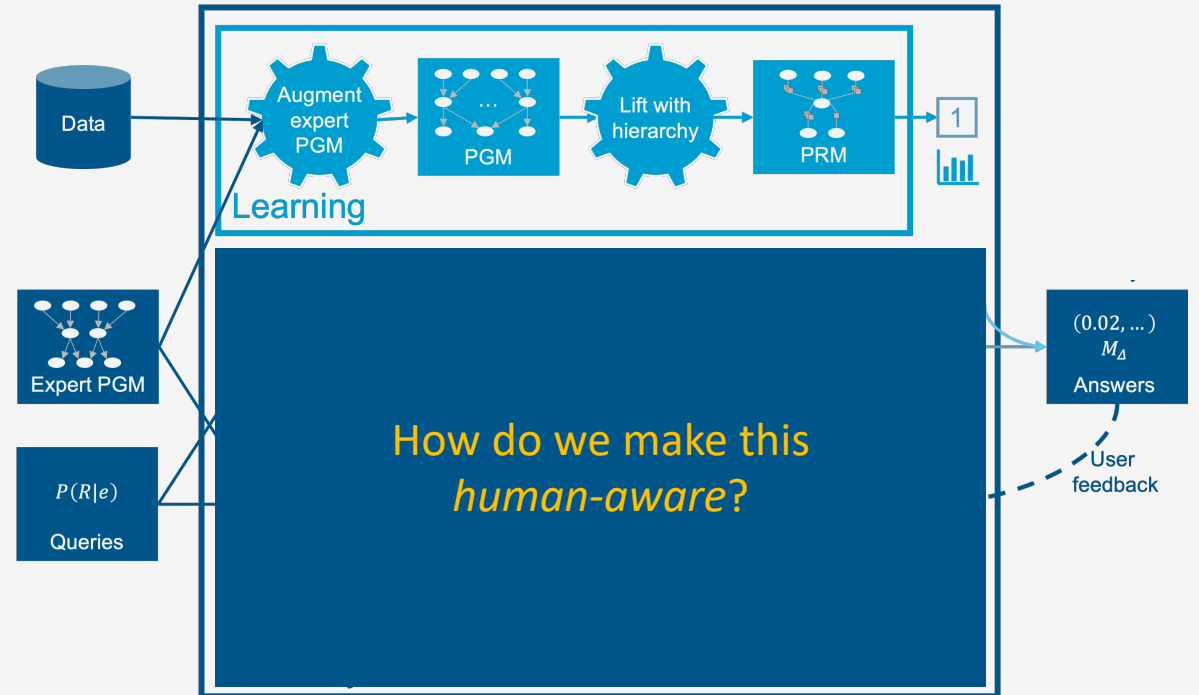
What If We Go Broader?

- Interpretable PGMs
 - All variables are observable (unrealistic!)
 - Assumption: Knowing the state of the system makes system explainable
 - Problem: Large nets are not readable even if state known
- Abstraction
 - *Lifted* versions for relational / symmetric domains:
 - Logical variables compactly encode repeated structures
 - Semantics: ground, multiply, normalise
 - Inference task defined as before, e.g.,
 $P(\text{Epid} | \text{Sick}(X) = \text{true})_{X \in \{\text{alice}, \text{eve}, \text{bob}\}}$
 - Assumption: Large nets become smaller and thus more easily readable



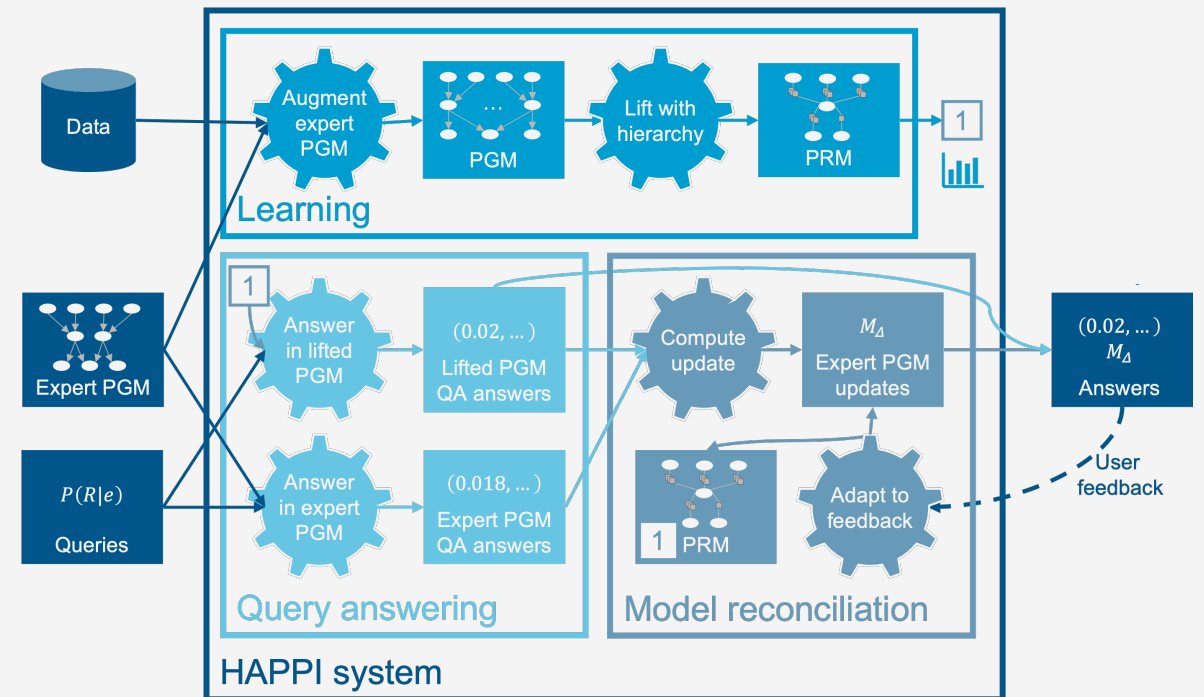
Human-aware Setting of Probabilistic Inference

- Models
 - Domain expert has domain knowledge
 - Implicit human's model
 - Domain expert builds a (small) PGM
 - Agent's model of human's model
 - Learning algorithm expands the PGM using additional training data
 - Can be much larger than the input PGM
 - Agent's model
- Inference
 - Human asks queries
 - Agent answers query



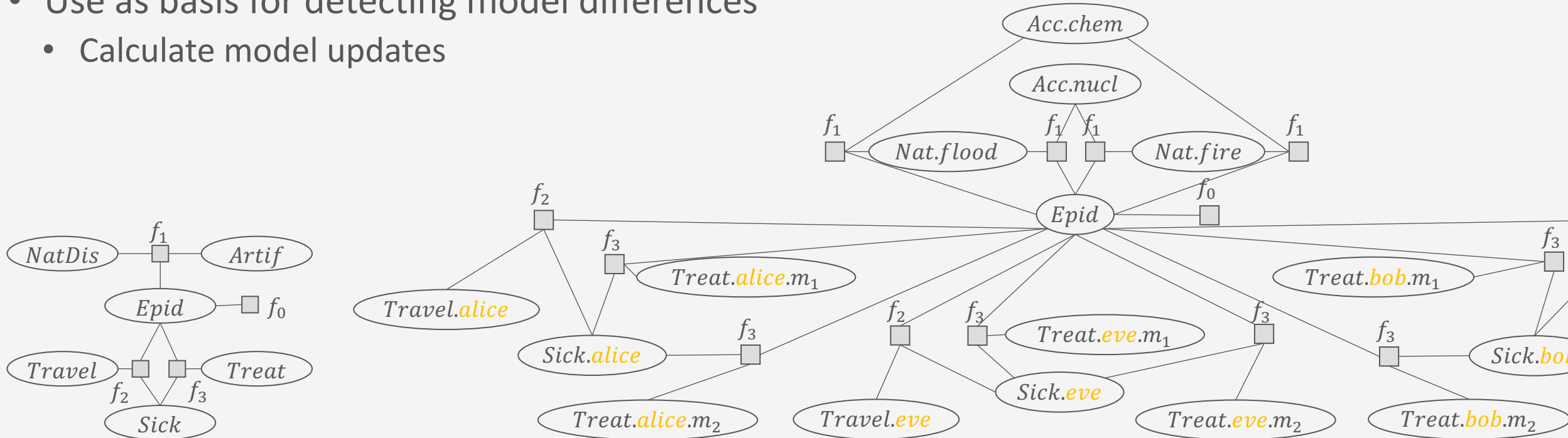
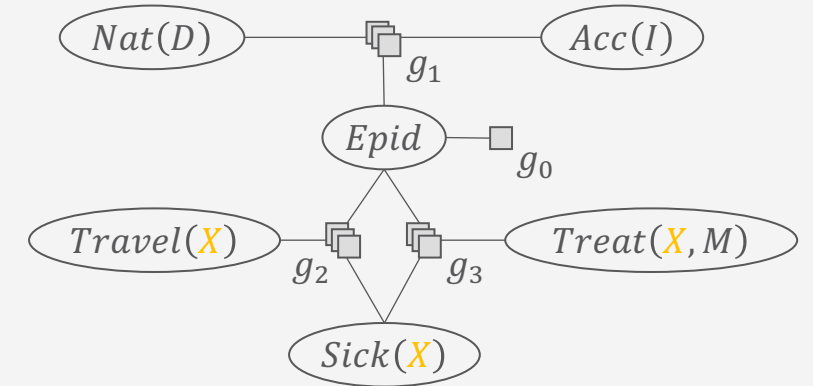
Model Reconciliation for Probabilistic Inference

- Reconciling model differences between expert model and trained model
- Use lifting to increasingly abstract large PGMs to get closer to smaller human's models
 - Hierarchy of models to be able to explain different parts in different detail
- Reconcile human's and agent's model by exposing differences to human
 - Expose only those parts relevant to a query to update human's model



Lifting as Abstraction

- Find (almost) repeated structures and abstract them
 - Colour passing as a variant of Weisfeiler-Leman
- Use as basis for detecting model differences
 - Calculate model updates



KI Starter: HAPPI Reconciliation

Human-aware PGMs and Probabilistic Inference via Lifted Model Reconciliation

- Goal: Provide explanations for probabilistic inference based on model reconciliation between trained and expert models
- Novel research direction that combines techniques that have exhibited great success in their respective areas
- Helps push fundamental task of inference into human-oriented AI

