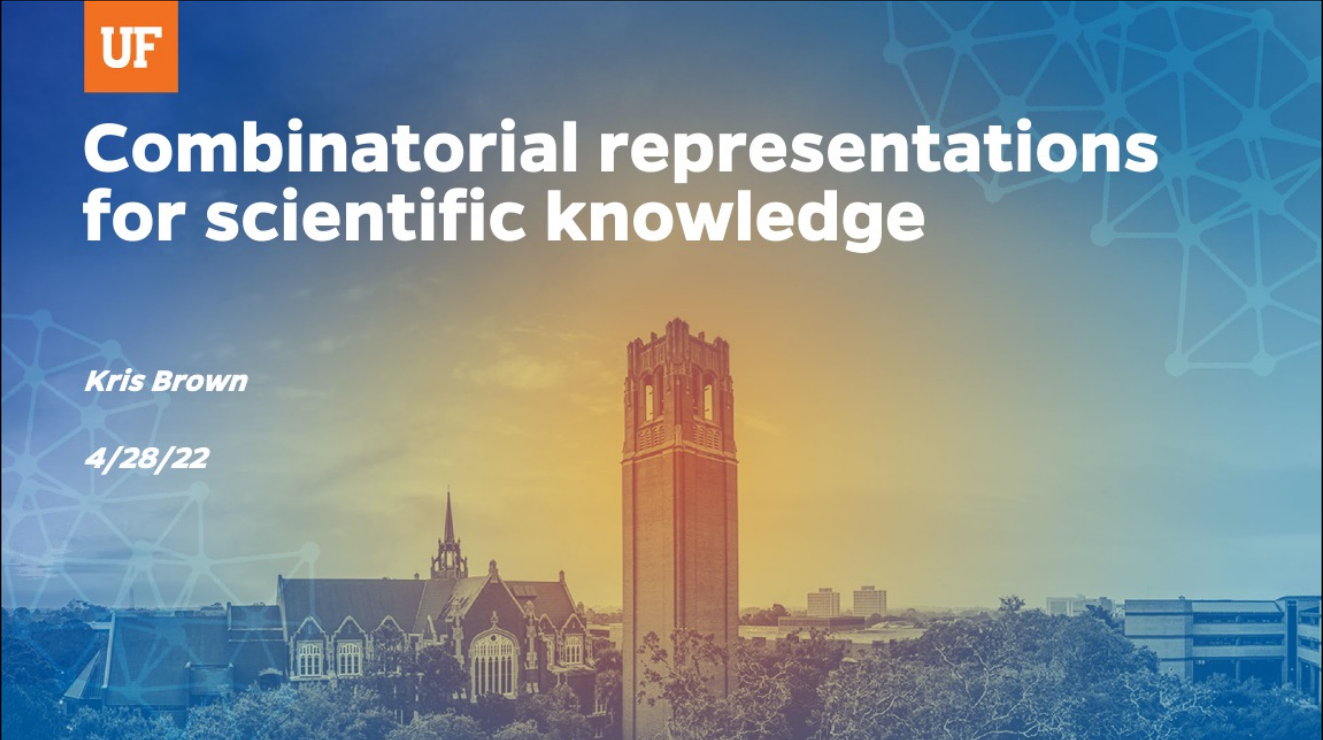


The logo for the University of Florida, consisting of the letters 'UF' in white on an orange square background.

Combinatorial representations for scientific knowledge

Kris Brown

4/28/22

A photograph of the University of Florida campus, featuring the prominent Spire tower and the Old Chapel building, set against a blue sky with a network diagram overlay. The image is overlaid with a blue gradient and a white network diagram consisting of nodes and connecting lines.

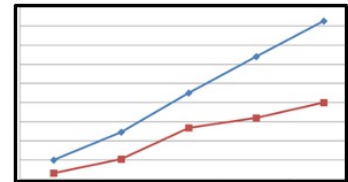
Hi, I'm a postdoc working with James Fairbanks for the past 8 months. Prior to that I was working as a computational chemist where I struggled a lot with issues that I think we're making great progress towards addressing. I hope if you're a computational scientist that this talk convinces you that these techniques are worth looking more into. This will be an applied category theory talk, so I'll focus on the applications rather than the theory.



Status quo: Model as opaque code



```
"""
2 H2 + O2 → 2 H2O. Mass-action kinetics. Compare to experimental data and plot.
"""
def main():
    # experimental data
    real_data = [0.0101, 0.012, 0.023, 0.037, 0.045, 0.053, 0.061, 0.069,
0.076, 0.083, 0.089, 0.096, 0.102, 0.108, 0.114, 0.119, 0.125, 0.130, 0.135,
0.140, 0.145, 0.150, 0.154, 0.159, 0.163, 0.167, 0.171, 0.175, 0.179, 0.183,
0.186, 0.190, 0.193, 0.197, 0.200, 0.203, 0.206, 0.209, 0.212, 0.215, 0.218,
0.221, 0.224, 0.227, 0.229, 0.232, 0.234, 0.237, 0.239]
    # Initial concentrations
    H2,O2, H2O = 1.0, 2.0, 0.0
    dt = 0.01
    results = []
    for step in range(1,50):
        print("Step ", step)
        rate = 0.5 * H2**2 * O2
        H2 -= 2*rate*dt
        O2 -= rate*dt
        H2O += rate*dt
        results.append(H2O)
    plot(results, real_data) # Figure 4 in the paper
```



- Knowledge informally represented in publications.
- Often the only formal representation is code.

2

I want to start by providing contrast with a very pervasive form of scientific knowledge, which seemed almost necessary to me while doing computational research.

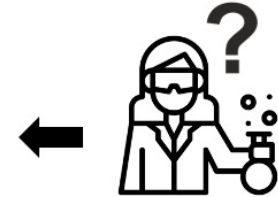
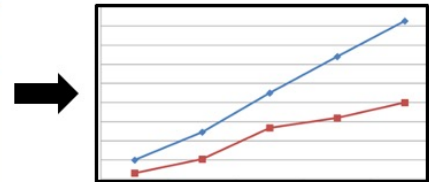
This is a small Python script which takes some data and compares it to a small simulation of a chemical reaction. This one is self-contained, but in general they live in a filesystem and can refer to each other. The big idea behind such a script is sometimes mentioned informally in a publication, and that is an important form of scientific knowledge but one we'll ignore for this talk because it is informal and can't be manipulated by tools (short of NLP).



Status quo: Model as opaque code



```
"""
2 H2 + O2 → 2 H2O. Mass-action kinetics. Compare to experimental data and plot.
"""
def main():
    # experimental data
    real_data = [0.0101, 0.012, 0.023, 0.037, 0.045, 0.053, 0.061, 0.069,
0.076, 0.083, 0.089, 0.096, 0.102, 0.108, 0.114, 0.119, 0.125, 0.130, 0.135,
0.140, 0.145, 0.150, 0.154, 0.159, 0.163, 0.167, 0.171, 0.175, 0.179, 0.183,
0.186, 0.190, 0.193, 0.197, 0.200, 0.203, 0.206, 0.209, 0.212, 0.215, 0.218,
0.221, 0.224, 0.227, 0.229, 0.232, 0.234, 0.237, 0.239]
    # Initial concentrations
    H2,O2, H2O = 1.0, 2.0, 0.0
    dt = 0.01
    results = []
    for step in range(1,50):
        print("Step ", step)
        rate = 0.5 * H2**2 * O2
        H2 -= 2*rate*dt
        O2 -= rate*dt
        H2O += rate*dt
        results.append(H2O)
    plot(results, real_data) # Figure 4 in the paper
```



- Knowledge informally represented in publications.
- Often the only formal representation is code.

3

If you point to a figure in the paper and ask the author "What does this mean?" and demand a rigorously defined object, you'll get pointed to this piece of code (though you'd also need to know lots of details about their computer environment: versions of libraries, environment variables, etc.). You could ask what are the downsides of this being the reality for many scientists.



Status quo: Model as opaque code



```
"""
2 H2 + O2 → 2 H2O. Mass-action kinetics. Compare to experimental data and plot.
"""
def main():
    # experimental data
    real_data = [0.0101, 0.012, 0.023, 0.037, 0.045, 0.053, 0.061, 0.069,
0.076, 0.083, 0.089, 0.096, 0.102, 0.108, 0.114, 0.119, 0.125, 0.130, 0.135,
0.140, 0.145, 0.150, 0.154, 0.159, 0.163, 0.167, 0.171, 0.175, 0.179, 0.183,
0.186, 0.190, 0.193, 0.197, 0.200, 0.203, 0.206, 0.209, 0.212, 0.215, 0.218,
0.221, 0.224, 0.227, 0.229, 0.232, 0.234, 0.237, 0.239]
    # Initial concentrations
    H2, O2, H2O = 1.0, 2.0, 0.0
    dt = 0.01
    results = []
    for step in range(1, 50):
        print("Step ", step)
        rate = 0.5 * H2**2 * O2
        H2 -= 2*rate*dt
        O2 -= rate*dt
        H2O += rate*dt
        results.append(H2O)
    plot(results, real_data) # Figure 4 in the paper
```

- Many tasks we'd like to do that cannot be done with arbitrary code (nor mathematical expressions).

The abstract critique of this style is that it mixes together data, structure, and semantics, which all should be kept apart.

The corresponding practical critique is lose the opportunity to do some cool things with our models if we make this mistake.



Status quo: Model as opaque code



```
"""
2 H2 + O2 → 2 H2O. Mass-action kinetics. Compare to experimental data and plot.
"""
def main():
    # experimental data
    real_data = [0.0101, 0.012, 0.023, 0.037, 0.045, 0.053, 0.061, 0.069,
0.076, 0.083, 0.089, 0.096, 0.102, 0.108, 0.114, 0.119, 0.125, 0.130, 0.135,
0.140, 0.145, 0.150, 0.154, 0.159, 0.163, 0.167, 0.171, 0.175, 0.179, 0.183,
0.186, 0.190, 0.193, 0.197, 0.200, 0.203, 0.206, 0.209, 0.212, 0.215, 0.218,
0.221, 0.224, 0.227, 0.229, 0.232, 0.234, 0.237, 0.239]
    # Initial concentrations
    H2, O2, H2O = 1.0, 2.0, 0.0
    dt = 0.01
    results = []
    for step in range(1,50):
        print("Step ", step)
        rate = 0.5 * H2**2 * O2
        H2 -= 2*rate*dt
        O2 -= rate*dt
        H2O += rate*dt
        results.append(H2O)
    plot(results, real_data) # Figure 4 in the paper
```

- Update code when assumptions change

• Many tasks we'd like to do that cannot be done with arbitrary code (nor mathematical expressions).

We'd like to have a language for declaring our assumptions and how they've changed, and when we do so, it's as if our code gets updated into the new framework.



Status quo: Model as opaque code



```
"""
2 H2 + O2 → 2 H2O. Mass-action kinetics. Compare to experimental data and plot.
"""
def main():
    # experimental data
    real_data = [0.0101, 0.012, 0.023, 0.037, 0.045, 0.053, 0.061, 0.069,
0.076, 0.083, 0.089, 0.096, 0.102, 0.108, 0.114, 0.119, 0.125, 0.130, 0.135,
0.140, 0.145, 0.150, 0.154, 0.159, 0.163, 0.167, 0.171, 0.175, 0.179, 0.183,
0.186, 0.190, 0.193, 0.197, 0.200, 0.203, 0.206, 0.209, 0.212, 0.215, 0.218,
0.221, 0.224, 0.227, 0.229, 0.232, 0.234, 0.237, 0.239]
    # Initial concentrations
    H2,O2, H2O = 1.0, 2.0, 0.0
    dt = 0.01
    results = []
    for step in range(1,50):
        print("Step ", step)
        rate = 0.5 * H2**2 * O2
        H2 -= 2*rate*dt
        O2 -= rate*dt
        H2O += rate*dt
        results.append(H2O)
    plot(results, real_data) # Figure 4 in the paper
```

- Update code when assumptions change
- Explore alternate reaction networks to fit the data

• Many tasks we'd like to do that cannot be done with arbitrary code (nor mathematical expressions).

There's a process that we all do of tweaking parameters and testing hypotheses until we get something that matches our experimental result – this process isn't made explicit or at all automated with this paradigm, even though there are certain mechanical and tedious aspects of it that beg to be automated.



Status quo: Model as opaque code



```
"""
2 H2 + O2 → 2 H2O. Mass-action kinetics. Compare to experimental data and plot.
"""
def main():
    # experimental data
    real_data = [0.0101, 0.012, 0.023, 0.037, 0.045, 0.053, 0.061, 0.069,
0.076, 0.083, 0.089, 0.096, 0.102, 0.108, 0.114, 0.119, 0.125, 0.130, 0.135,
0.140, 0.145, 0.150, 0.154, 0.159, 0.163, 0.167, 0.171, 0.175, 0.179, 0.183,
0.186, 0.190, 0.193, 0.197, 0.200, 0.203, 0.206, 0.209, 0.212, 0.215, 0.218,
0.221, 0.224, 0.227, 0.229, 0.232, 0.234, 0.237, 0.239]
    # Initial concentrations
    H2,O2, H2O = 1.0, 2.0, 0.0
    dt = 0.01
    results = []
    for step in range(1,50):
        print("Step ", step)
        rate = 0.5 * H2**2 * O2
        H2 -= 2*rate*dt
        O2 -= rate*dt
        H2O += rate*dt
        results.append(H2O)
    plot(results, real_data) # Figure 4 in the paper
```

- Update code when assumptions change
- Explore alternate reaction networks to fit the data
- Generate the entire code from just declaring the reaction

• Many tasks we'd like to do that cannot be done with arbitrary code (nor mathematical expressions).

In some sense writing the code at all seems like an undesirable, error-prone task. If you look at the description at the top of the function, we ought to be able to just declare that information and get the whole simulator for free.



Status quo: Model as opaque code



```
"""
2 H2 + O2 → 2 H2O. Mass-action kinetics. Compare to experimental data and plot.
"""
def main():
    # experimental data
    real_data = [0.0101, 0.012, 0.023, 0.037, 0.045, 0.053, 0.061, 0.069,
0.076, 0.083, 0.089, 0.096, 0.102, 0.108, 0.114, 0.119, 0.125, 0.130, 0.135,
0.140, 0.145, 0.150, 0.154, 0.159, 0.163, 0.167, 0.171, 0.175, 0.179, 0.183,
0.186, 0.190, 0.193, 0.197, 0.200, 0.203, 0.206, 0.209, 0.212, 0.215, 0.218,
0.221, 0.224, 0.227, 0.229, 0.232, 0.234, 0.237, 0.239]
    # Initial concentrations
    H2,O2, H2O = 1.0, 2.0, 0.0
    dt = 0.01
    results = []
    for step in range(1,50):
        print("Step ", step)
        rate = 0.5 * H2**2 * O2
        H2 -= 2*rate*dt
        O2 -= rate*dt
        H2O += rate*dt
        results.append(H2O)
    plot(results, real_data) # Figure 4 in the paper
```

- Update code when assumptions change
- Explore alternate reaction networks to fit the data
- Generate the entire code from just declaring the reaction
- Check if another model is the same / a submodel

• Many tasks we'd like to do that cannot be done with arbitrary code (nor mathematical expressions).

We have many informal notions of when models are equivalent or submodels of each other in some sense, and regardless of what we pick we cannot actually compute that without a better formal representation of the model.



Status quo: Model as opaque code



```
"""
2 H2 + O2 → 2 H2O. Mass-action kinetics. Compare to experimental data and plot.
"""
def main():
    # experimental data
    real_data = [0.0101, 0.012, 0.023, 0.037, 0.045, 0.053, 0.061, 0.069,
0.076, 0.083, 0.089, 0.096, 0.102, 0.108, 0.114, 0.119, 0.125, 0.130, 0.135,
0.140, 0.145, 0.150, 0.154, 0.159, 0.163, 0.167, 0.171, 0.175, 0.179, 0.183,
0.186, 0.190, 0.193, 0.197, 0.200, 0.203, 0.206, 0.209, 0.212, 0.215, 0.218,
0.221, 0.224, 0.227, 0.229, 0.232, 0.234, 0.237, 0.239]
    # Initial concentrations
    H2,O2, H2O = 1.0, 2.0, 0.0
    dt = 0.01
    results = []
    for step in range(1,50):
        print("Step ", step)
        rate = 0.5 * H2**2 * O2
        H2 -= 2*rate*dt
        O2 -= rate*dt
        H2O += rate*dt
        results.append(H2O)
    plot(results, real_data) # Figure 4 in the paper
```

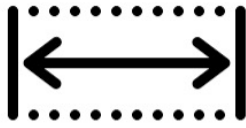
- Update code when assumptions change
- Explore alternate reaction networks to fit the data
- Generate the entire code from just declaring the reaction
- Check if another model is the same / a submodel
- Easily alter semantics (e.g. stochastic-based simulation)

• Many tasks we'd like to do that cannot be done with arbitrary code (nor mathematical expressions).

Lastly, once you make some syntax-semantics distinction, it becomes clear that we should be able to give the same syntax a different semantics, such as taking this chemical reaction and instead running a stochastic individual-based simulation rather than the analytical, aggregate ODE simulation that you see here.

Overall, there are many tasks that we'd like to do that can't be done with arbitrary code. We'll see even giving a thorough specification of our system in the language of differential equations does not help us towards these goals.

Outline



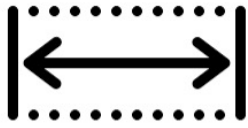
Combinatorial
Data Structures



Epidemiology
Case Studies

So I want to this talk in two parts, at first giving a bird's eye view with a lot of breadth, and then picking a particular type of scientific model and going in depth.

Outline



Combinatorial
Data Structures



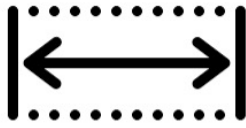
Epidemiology
Case Studies



Applications: Where?

So I'll first say where are real world domains in which we'll apply combinatorial structures

Outline



Combinatorial
Data Structures



Epidemiology
Case Studies



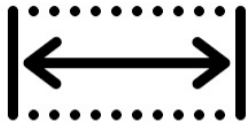
Applications: Where?



Structures: What?

Then I'll get more into examples of what they are.

Outline



Combinatorial
Data Structures



Epidemiology
Case Studies



Applications: Where?



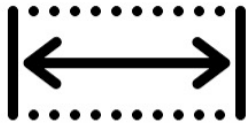
Structures: What?



Algorithms: How?

I'll highlight some algorithms that can operate on them.

Outline



Combinatorial
Data Structures



Epidemiology
Case Studies



Applications: Where?



Structures: What?



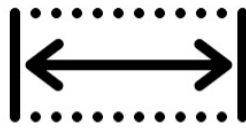
Algorithms: How?



Virtues: Why?

And then try to summarize what are the virtues of using them over some standard alternatives.

Outline



Combinatorial Data Structures



Applications: Where?



Structures: What?



Algorithms: How?



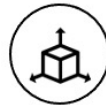
Virtues: Why?



Epidemiology Case Studies



Specification



Exploration



Simulation

The second part will focus on Petri Net models in epidemiology and talk about declaring, exploring, and simulating them all in a transparent way.

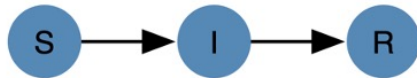
First half I'll make a lot of claims and you'll have to take my word for it... But in the second half I'll give more evidence for the claims I make... though a rigorous understanding of why this works requires some category theory.



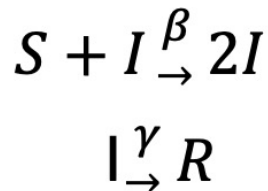
Where: model specification



How to represent the data of a scientific model?



(informal diagram)



(formal expression tree)

So to begin with this question, note that a standard solution is to use some sort of informal diagram along with natural language to convey intuition, and then provide some mathematical equations or code which are unambiguous and precise.

As an example which will be featured many times in this presentation, here is a very simple kind of model used in epidemiology called the SIR model. The intuition is that people start out susceptible to some disease, such as COVID. They then can become infected with some probability, and then they transition to recovered with some probability. The expressions that make this precise actually say that there exist two transition events: one susceptible person and one infected person combine to make two infected people, and one infected person becomes one recovered person.



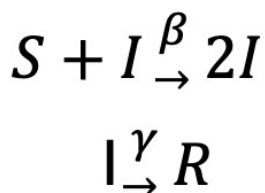
Where: model specification



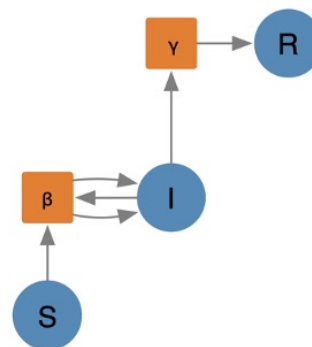
How to represent the data of a scientific model?



(informal diagram)



(symbolic expressions)



(Petri Net – a kind of formal diagram)

Petri nets: a rigorous, diagrammatic syntax for the *structure* of chemical reaction networks.

An improvement that we would propose is that there exists a type of diagram which is well defined and intuitive like a graph, but it actually can be unambiguously converted into those formal reaction expressions. This is called a Petri net, and you can see how the multiple arrows are keeping track of the data in the formal expression. I will give a formal definition for it later in the talk.

That might seem like a superficial difference from the equation formalism, but consider now the addition of constraints to some class of reactions we want to consider.



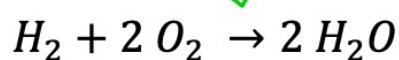
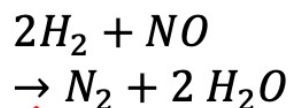
Where: model specification



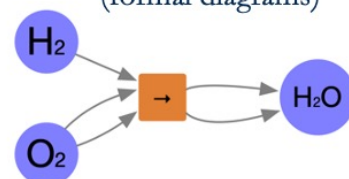
How to represent constraints associated with chemical reaction networks?

- No more than two reactants per reaction

(formal expression trees)



(formal diagrams)



Standard syntax for reaction networks can't enforce constraints *by construction*, but a variant of Petri Nets can. More details in Part 2 of the talk.

In the case of chemistry, we often want to exclude from consideration reactions that contain three things that react at the same time. And you can ask: What data structure do you represent your reactions in such that you automatically enforce this constraint?



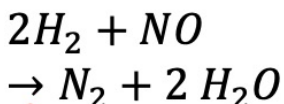
Where: model specification



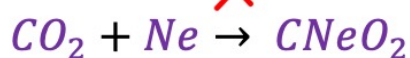
How to represent constraints associated with gas-phase chemistry?

- No more than two reactants per reaction
- All species are either **Reactive** or **Inert** (Inert species do not react)

(formal expression trees)



X



X

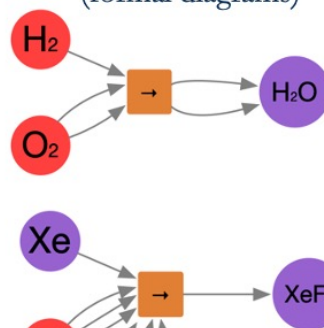


✓



✓

(formal diagrams)



Standard syntax for reaction networks can't enforce constraints *by construction*, but a variant of Petri Nets can. More details in Part 2 of the talk.

I give another more sophisticated kind of constraint below, saying that all species must be tagged as either inert or reactive and furthermore you can't have a reaction between two inert species.

The punchline will come in the second half of the talk, where I describe something very closely related to the petri net you see here that IS able to represent these constraints purely by construction, meaning all possible diagrams you can draw correspond precisely to the subset of reactions which satisfy these logical constraints.

So this is something you can do working in this alternative syntax to the syntax of mathematical expression trees. It's less powerful in what it can represent, but more powerful in what we can do with it. (that's a recurring theme)



Where: model comparison



- When are two models structurally the same, in some sense?
- When are two models behaviorally the same, in some sense?

$$\left(\begin{array}{ccc} C : \Omega_t^0 & \xrightarrow{\partial_t} & \dot{C} : \Omega_t^0 \xleftarrow{\star^{-1}} d\phi : \tilde{\Omega}_t^3 \\ d \downarrow & & \uparrow d \\ dC : \Omega_t^1 & \xrightarrow{k\star} & \phi : \tilde{\Omega}_t^2 \end{array} \right) \rightarrow \left(C : \Omega_t^0 \xrightarrow[k\Delta]{\partial_t} \dot{C} : \Omega_t^0 \right)$$

(Two formal graphical presentations of heat diffusion)

Diagrammatic representation of multiphysics allows us to check behavioral equivalence, unlike working with the raw PDEs

Patterson et al. "A diagrammatic view of differential equations in physics" (2022)

Another application of combinatorial data structures is to compare whether two models are the same, structurally or behaviorally.

In the example here I'm showing, a different kind of formal, graphical representations for differential equations is used to represent the heat equations in two different ways (things can be named differently and structured differently).

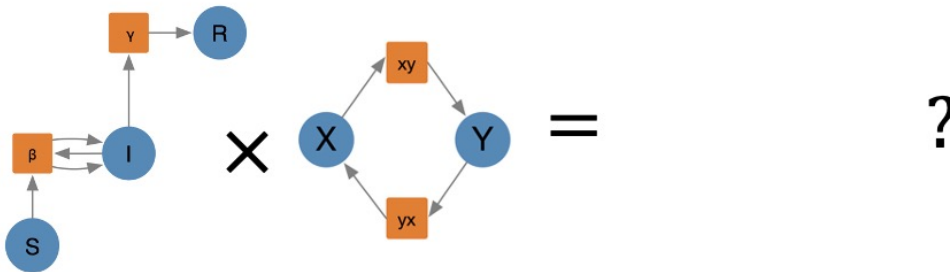
But in virtue of how they are represented, we have a decidable algorithm which computes behavioral similarity by determining whether or not the solutions of one system of equations are all solutions of the other. This is not something you can do if just handed a bunch of PDEs, because again raw mathematical expressions have way too much expressive power, so in general it's rare for decidable algorithms that do useful things with them.



Where: model exploration



- How to combine basic ideas for models into complex models?
- How to efficiently use experimental data to pick the best model?



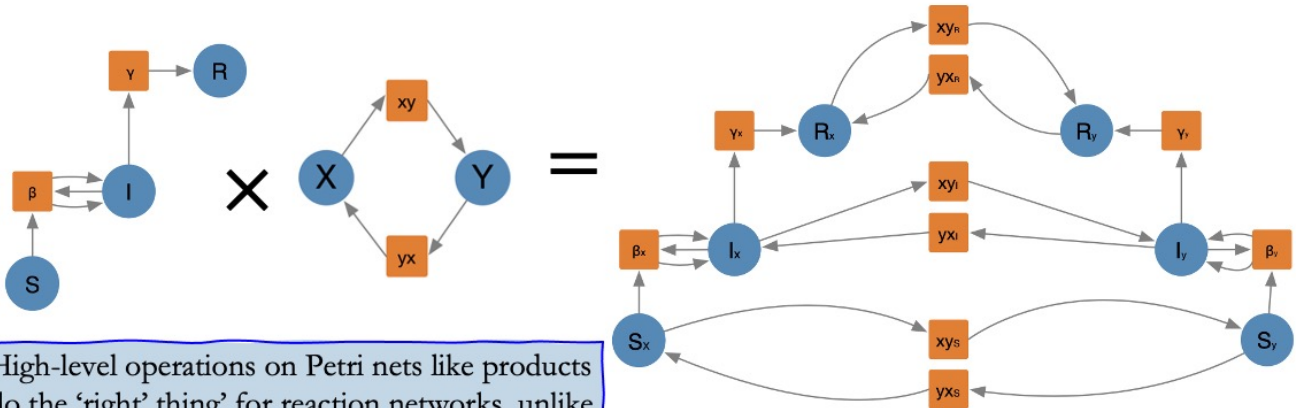
Another place we want to deploy these representations is the model exploration problem. I'll go more into detail about this later, but as a teaser, consider what it would mean to combine two Petri net models, one of which is the SIR model and another which we call a "two city model", where people live in United States and Europe and can move between them.



Where: model exploration



- How to combine basic ideas for models into complex models?
- How to efficiently use experimental data to pick the best model?



High-level operations on Petri nets like products do the 'right' thing' for reaction networks, unlike for symbolic syntax or raw ODEs.

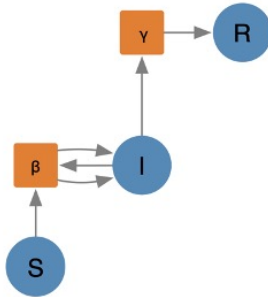
The intuitive answer can be elegantly using these Petri net representations (category theory tells us what the right notion of product is for them), whereas taking the product of the literal differential equations wouldn't yield anything useful.



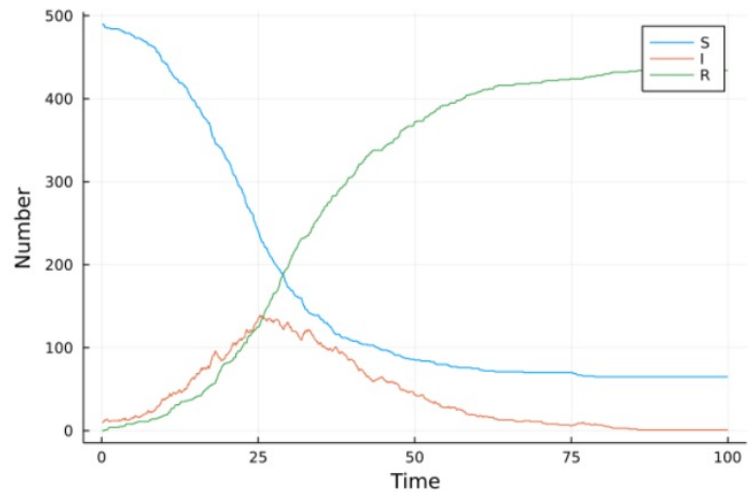
Where: model inference



How to apply a model to a concrete problem?

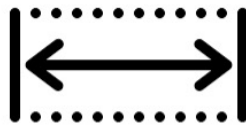


Specifying a reaction network as a Petri Net allows for automatically generating a simulator program.



The last application area I want to highlight is the problem of model inference. If you were handed this petri net on the left, could you automatically generate a simulator which accepts initial conditions and reaction rates and produces a plot on the right? This is a useful feature because scientists really shouldn't be writing code – we envision a future where they can work purely in their true domain, and with the right representations of their models the solvers should be automatically generated.

Outline



Combinatorial Data Structures



Applications: Where?



Structures: What?



Algorithms: How?



Virtues: Why?



Epidemiology Case Studies



Specification



Exploration



Simulation

Alright that was whirlwind of applications, but it's good to know the motivation for why we're doing this. Now I'll say more about what these combinatorial structures are.

Now the next slide has a lot of examples of the data structures that make those applications possible. CAVEAT: it will be abstract and you don't actually need to fully grasp them in order to use them or understand the applications. What I'll show is what makes the math behind the picture



What: combinatorial data structures



Datatype	Set
Structure	A
Visualized example	

Sets, Functions, and Graphs are all simple examples of combinatorial data structures.

So combinatorial data structures I think of as elements that are connected. Here are some basic examples. Starting with a set which is just a bunch of unconnected elements, no structure



What: combinatorial data structures



Datatype	Set	Function
Structure	A	$A \xrightarrow{f} B$
Visualized example		

Sets, Functions, and Graphs are all simple examples of combinatorial data structures.

Functions then allow us to connect elements living in different sets.



What: combinatorial data structures



Datatype	<u>Set</u>	<u>Function</u>	<u>Graph</u>
Structure	A	$A \xrightarrow{f} B$	$E \begin{matrix} \xrightarrow{\text{src}} \\ \xrightarrow{\text{tgt}} \end{matrix} V$
Visualized example			

Sets, Functions, and Graphs are all simple examples of combinatorial data structures.

Directed graphs are really just a pair of functions, where the edges themselves are a set and the functions say where their source and target are.

As you can see, each of these examples is going up a level of abstraction. The next level of abstraction requires us to think of the structure itself as a variable parameter.



What: combinatorial data structures



Datatype	Set	Function	Graph
Structure	A	$A \xrightarrow{f} B$	$E \begin{matrix} \xrightarrow{\text{src}} V \\ \xrightarrow{\text{tgt}} V \end{matrix}$
Visualized example			



These are C-Sets, too.
The "C" is the Structure.

Datatype Class	C-Set
Example Structure 'C'	
Visualized example (with above structure)	

- C-Sets are a powerful abstraction.
- Each "C" gives a new datatype.
- Instance is network of sets+functions.
- C-Sets can represent anything you would put in a database
- C-Sets are both expressive and easy to compute with.

Call this parameter "C" (it turns out it's something called a category, but think of it as a graph). That's to say, we're now working with data structures, where the shape of the structure is itself a parameter. This is really powerful because it means if you write a C-set algorithm to work with arbitrary C, it's like you've implemented an algorithm for a massive variety of very different looking datatypes.



What: combinatorial data structures



Datatype	Set	Function	Graph
Structure	A	$A \xrightarrow{f} B$	$E \begin{matrix} \xrightarrow{\text{src}} \\ \xrightarrow{\text{tgt}} \end{matrix} V$
Visualized example			

Datatype Class	C-Set
Example Structure 'C'	

Visualized example (with above structure)	
--	--

C-Set

Slice

Span

Structured cospan

Diagram

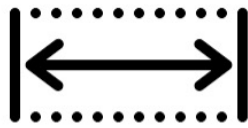
Sketch

Attributed C-Set

We'll use many other structures that build upon C-Sets.

On my list of structures on the right, C-sets are just the starting point, and all the ones below it actually build more abstractly upon C-sets. But we'll introduce those as needed once we start talking about the key projects I'd like to showcase.

Outline



Combinatorial
Data Structures



Epidemiology
Case Studies



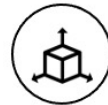
Applications: Where?



Specification



Structures: What?



Exploration



Algorithms: How?



Simulation

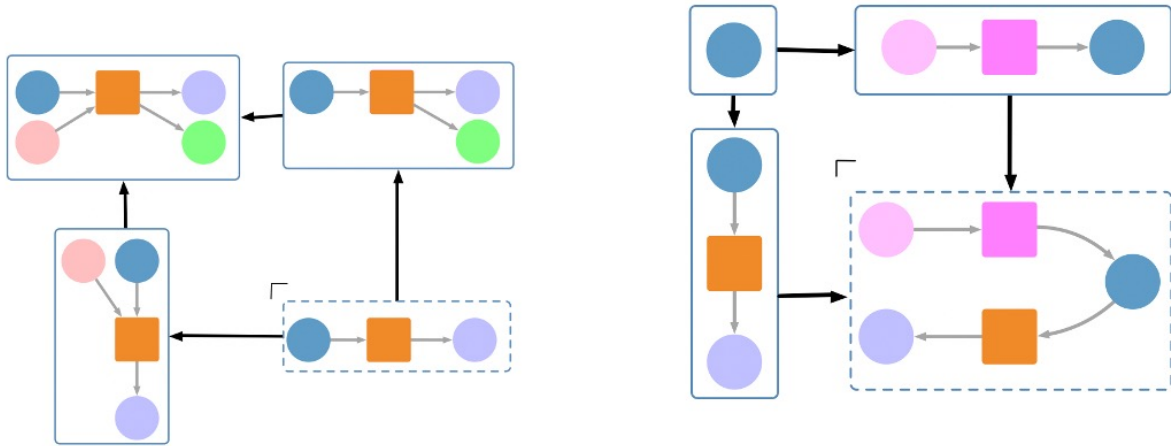


Virtues: Why?

I wanted to give a quick overview of some algorithms that operate on these combinatorial structures but that gets a bit technical and it's probably best to show applications of those algorithms, so I'll leave skip those slides for now.



How: Limits and Colimits



HIDDEN

Limits and colimits respectively, loosely correspond to the ideas of solving equations and gluing structures together.



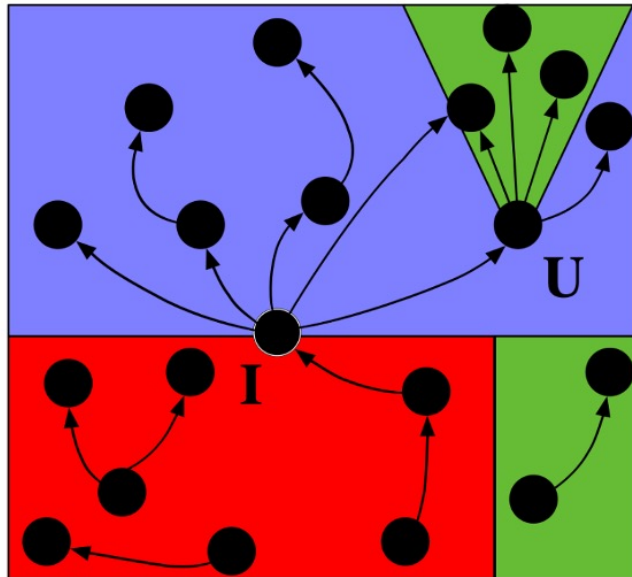
How: The Chase



Satisfies Σ

$I \rightarrow J$

$I \not\rightarrow J$

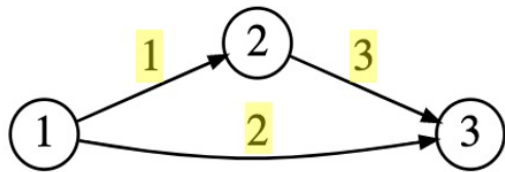


HIDDEN

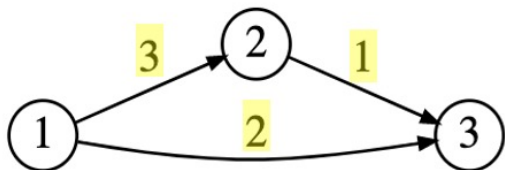
The chase is an algorithm which takes a model, call it "I" there at the center of this cartoon, and it tries to make it satisfy some constraints which we'll call sigma. It doesn't want just any model that satisfies this though. If you imagine these arrows as adding assumptions to a model, then the algorithm is finding the best possible way to make I satisfy the equations, by adding the least possible assumptions. This shows up in a surprising number of contexts.



How: Canonical isomorph of C-Sets





\cong

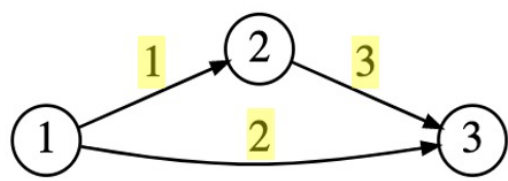


HIDDEN

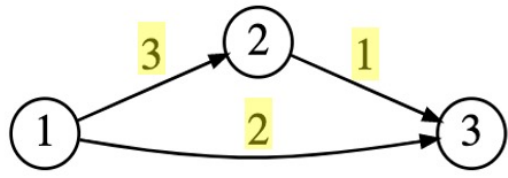
Lastly, there is this notion of symmetry where two models really are morally the same even though they've been labeled differently. This is isomorphism.

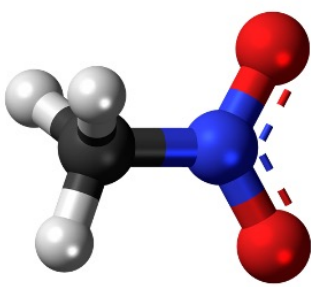
This algorithm finds a specific labelling that is equal for all isomorphic C-Sets.

 How: Canonical isomorph of C-Sets 





\cong








↓ ?


www.ccdc.cam.ac.uk


 COLLABORATIVE DRUG DISCOVERY



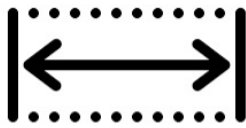



 ChemSpider
The free chemical database

HIDDEN

This is relevant for scientists who, for example, might have a molecule and want to check if it exists in some massive databases. If we don't compute a canonical representation of ordering all the atoms and bonds, then it can be very inefficient to check, for each molecule in all these databases, whether or not there exists an isomorphism.

Outline



Combinatorial Data Structures



Applications: Where?



Structures: What?



Algorithms: How?



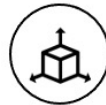
Virtues: Why?



Epidemiology Case Studies



Specification



Exploration



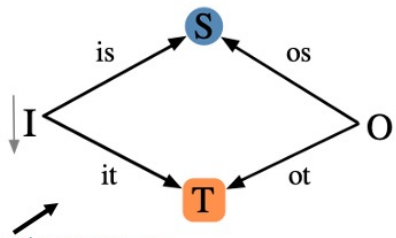
Simulation

I'd rather focus on why these combinatorial structures are a healthier way of representing knowledge.

The theme of these virtues is that there are really nice things you want to do with models but cannot do them if models are code or math. You can only do them when you represent your model as these (rigorous) pictures.



Why: Transparency / predictability



(model does not incorporate any concepts other than S-T connectivity)

- See what assumptions were made immediately
- Model is *data*, rather than code. Can be analyzed.

Analyzing the behavior of combinatorial data structures is easier than arbitrary code

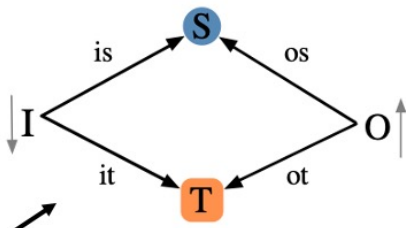
So what are these benefits? The first one I call transparency or predictability.

By separating the structure of the model from its behavior, you can see assumptions that have been made.

Without even looking at a particular Petri Net, by knowing the structure of it's schema we can see what sorts of concepts it could possibly contain.



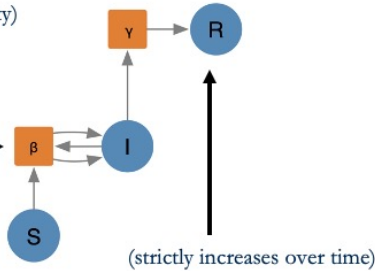
Why: Transparency / predictability



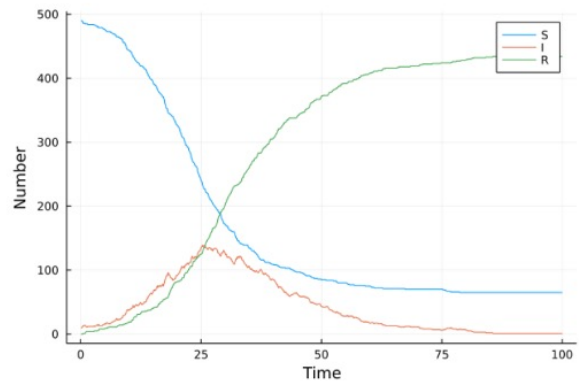
- See what assumptions were made immediately
- Model is *data*, rather than code. Can be analyzed.

(model does not incorporate any concepts other than S-T connectivity)

(rate does not depend on # of recovered people)



(strictly increases over time)



(for all possible simulations)

Analyzing the behavior of combinatorial data structures is easier than arbitrary code

Then, looking at the structure of the Petri net, even a simple computer algorithm can deduce that the number of recovered people will strictly increase for all possible simulations.

This is pretty simple stuff here, but when I talk later about encoding constraints into the structure, you'll see this is a powerful point.

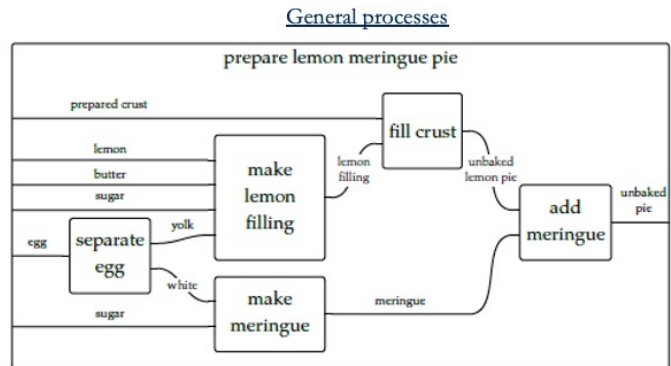
In general, analyzing behavior is easier than with arbitrary code or arbitrary math.



Why: Generality of scope



- Many different applications share the same syntax
- E.g. these are all applications of *wiring diagrams* (a particular kind of C-Set)



Explicit distinction between syntax and semantics makes it easy to generalize models

Fong and Spivak. "Seven Sketches in Compositionality" (2018)

So this slide visualizes a directed wiring diagram, which is actually a specific type of C-Set. You can see here it represents general processes such as making a pie.

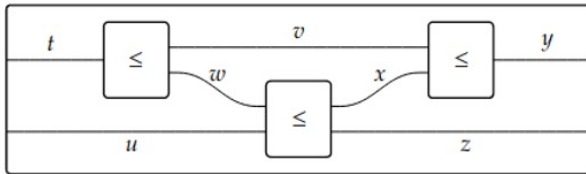


Why: Generality of scope



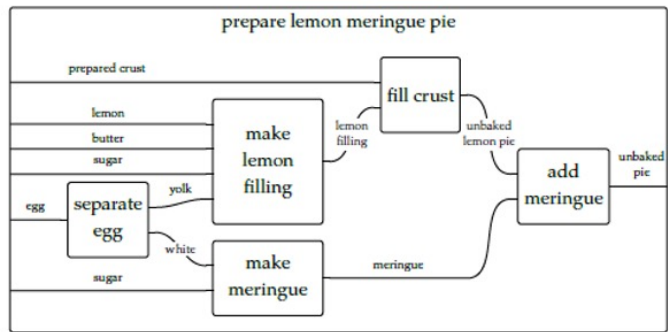
- Many different applications share the same syntax
- E.g. these are all applications of *wiring diagrams* (a particular kind of C-Set)

Inequality reasoning



$t \leq v + w$ and $w + u \leq x + z$ and $v + x \leq y$
 implies $t + u \leq y + z$

General processes



Explicit distinction between syntax and semantics makes it easy to generalize models

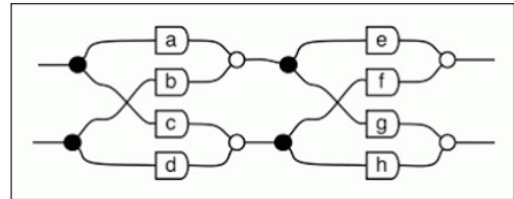
Fong and Spivak. "Seven Sketches in Compositionality" (2018)

But this same syntax allows us to represent mathematical proofs of a certain kind.



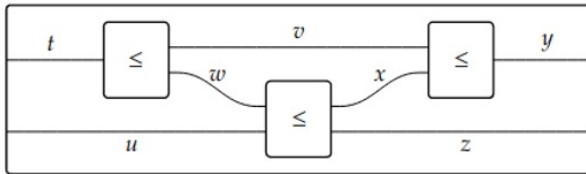
Why: Generality of scope

- Many different applications share the same syntax
- E.g. these are all applications of *wiring diagrams* (a particular kind of C-Set)



Pawel Sobocinski (<http://graphicallinearalgebra.net>)

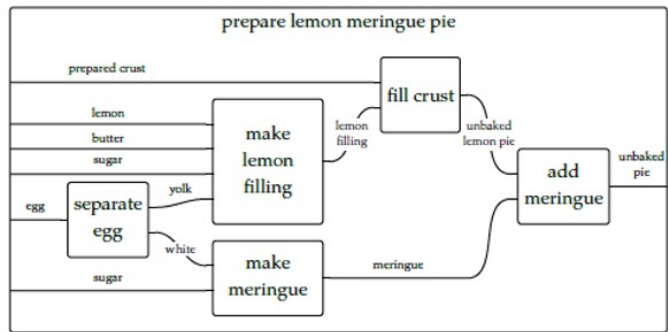
Inequality reasoning



$t \leq v + w$ and $w + u \leq x + z$ and $v + x \leq y$
 implies $t + u \leq y + z$

Explicit distinction between syntax and semantics makes it easy to generalize models

General processes



Fong and Spivak. "Seven Sketches in Compositionality" (2018)

And furthermore we can even represent linear algebra with these diagrams, which becomes a lot less mysterious when it is visual. (THERE'S A LEARNING CURVE.)

When you write code for wiring diagrams in general, it becomes applicable in a wide variety of circumstances (obviously, also code for C-Sets in general has wide applications).

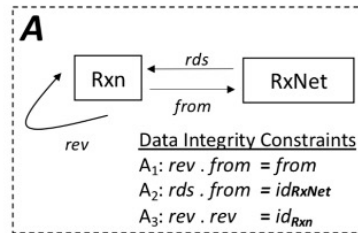
Basically, syntax and semantics are two independent knobs we often want to tweak, so keeping them separate makes it very easy to do the kinds of generalizations that we naturally want to do.



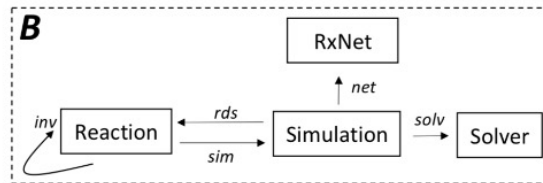
Why: Transferability



Migrate / merge data in different schemas



Data can be automatically moved between different structures when our models change. A tedious and error-prone process otherwise.



Brown, Spivak, Wisnesky. "Computational Data integration for Computational Science" (2019).

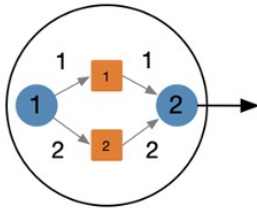
Generality is also related to transferability. If one's model of the world updates and you want to migrate your infrastructure from the old to the new, it's possible to do this knowing just from declaring the relationship of the structure of the old data to the new data. This can't be done when the model is a uniform block of arbitrary programming language code.



Why: Hierarchical design



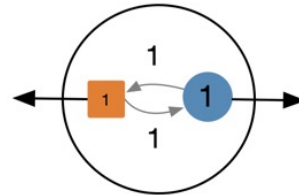
“Find all states which are the output of two transitions”



```

SELECT state2.id
FROM   S AS statel, S AS state2,
        T AS tran1,  T AS tran2,
        I AS in1,   I AS in2,
        O AS out1,  O AS out2
WHERE in1.is = statel.id,
        in1.it = tran1.id
        out1.os = state2.id
        out1.ot = tran1.id
        ...
  
```

“Find all catalysts and the reactions they catalyze”



```

SELECT tran1.id, statel.id
FROM   S AS statel, T AS tran1,
        I AS in1,   O AS out1
WHERE in1.is = statel.id,
        in1.it = tran1.id
        out1.os = statel.id
        out1.ot = tran1.id
  
```

I want to highlight the hierarchical design virtue by considering querying a database. If you're not familiar with this problem, just imagine that, if I have a gigantic petri net, we need to use a special language that can efficiently extract data from it. This language, SQL, is at the bottom here.

But we also have a special graphical syntax which corresponds to SQL queries, and this representation has an important advantage.

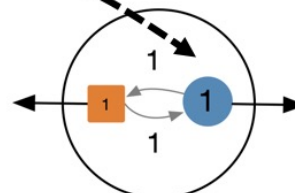
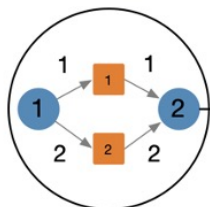


Why: Hierarchical design



Database queries can be built hierarchically using a wiring diagram syntax. Raw SQL queries are not composable this way.

“Find all catalysts (*that are the product of two reactions*) and the reactions they catalyze”



```

SELECT state2.id
FROM   S AS state1, S AS state2,
       T AS tran1,  T AS tran2,
       I AS in1,   I AS in2,
       O AS out1,  O AS out2
WHERE  in1.is = state1.id,
       in1.it = tran1.id
       out1.os = state2.id
       out1.ot = tran1.id
...

```

```

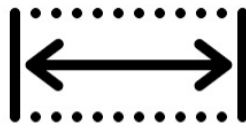
SELECT tran1.id, state1.id
FROM   S AS state1, T AS tran1,
       I AS in1,   O AS out1
WHERE  in1.is = state1.id,
       in1.it = tran1.id
       out1.os = state1.id
       out1.ot = tran1.id

```

The picture form has a special advantage in that it is compositional – you could take the query on the left and substitute it into the query on the right.

SQL code in contrast, doesn't let you do this kind of substitution (you'd have to write a special purpose algorithm to handle all the edge cases)

Outline



Combinatorial Data Structures



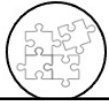
Applications: Where?



Structures: What?



Algorithms: How?



Virtues: Why?



Epidemiology Case Studies



Specification



Exploration

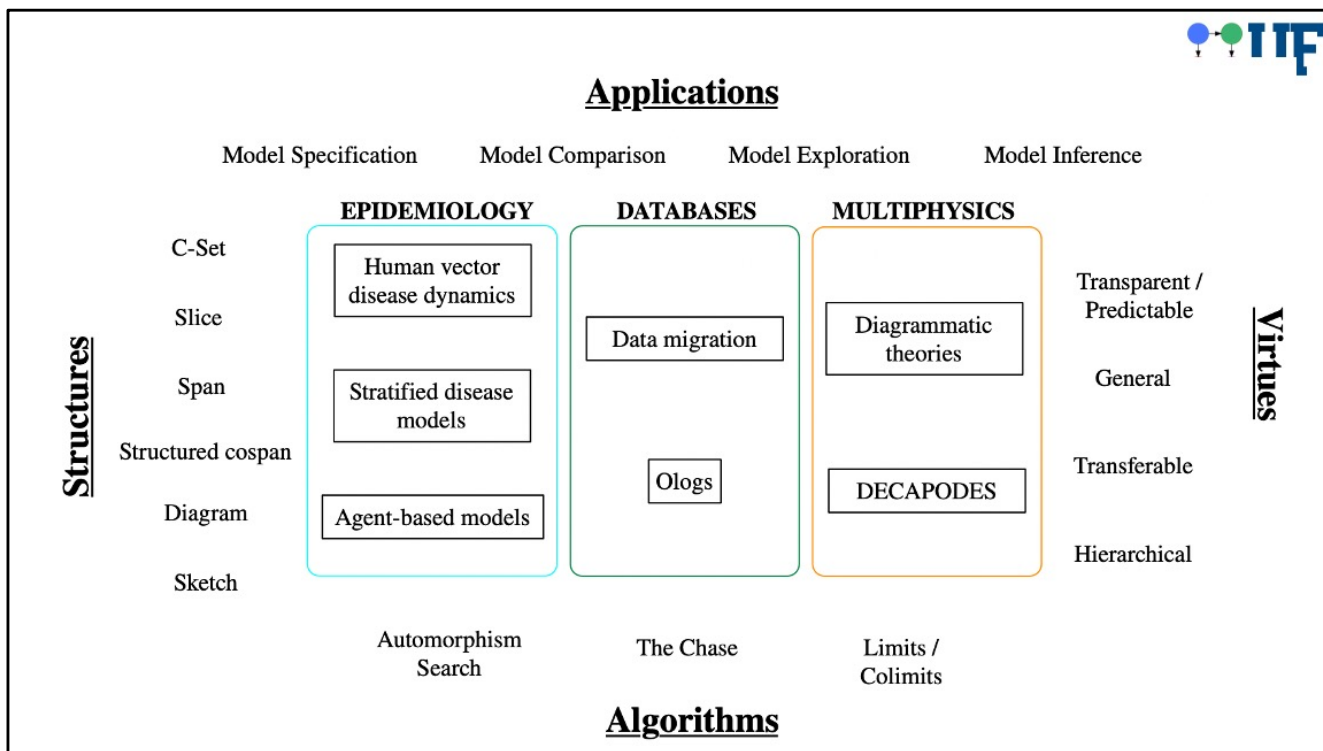


Simulation

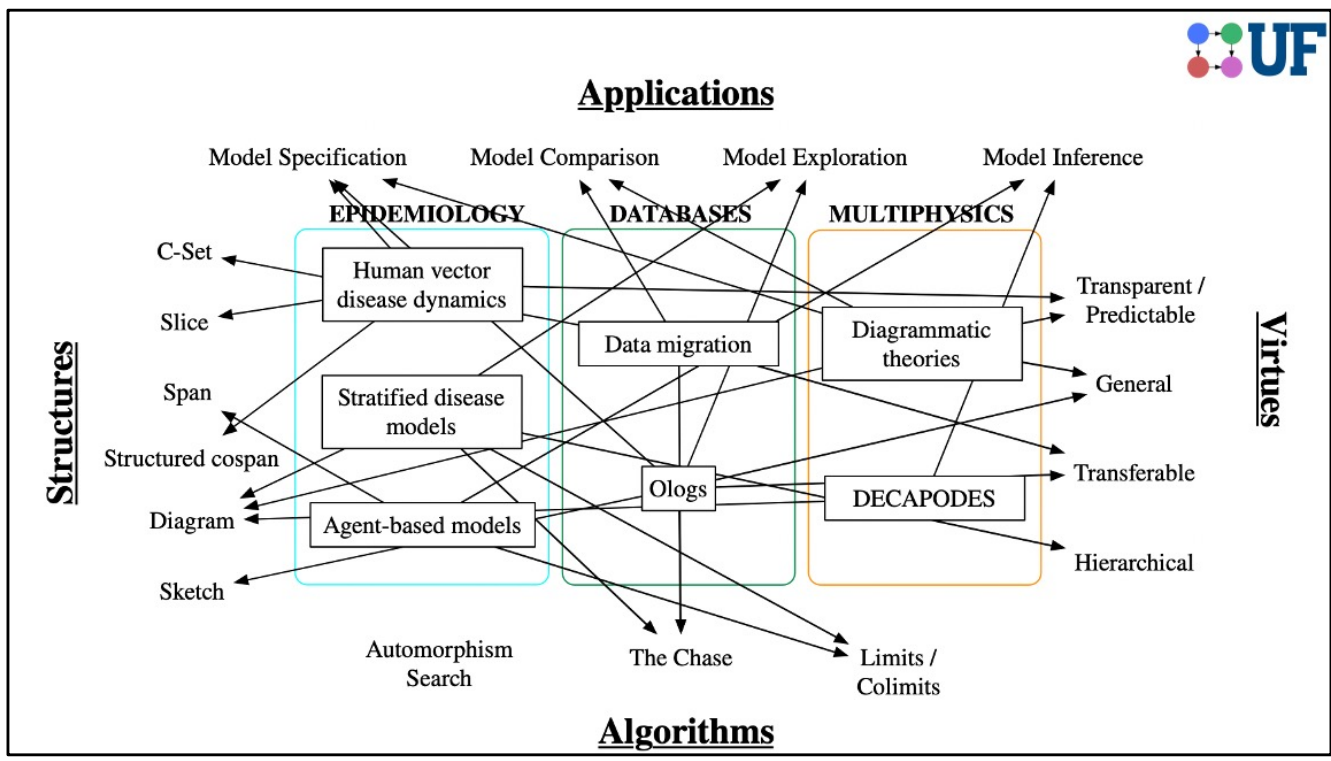
So that was very broad and very fast.

RECAP: lots of things you can't do with models as code that we highlighted

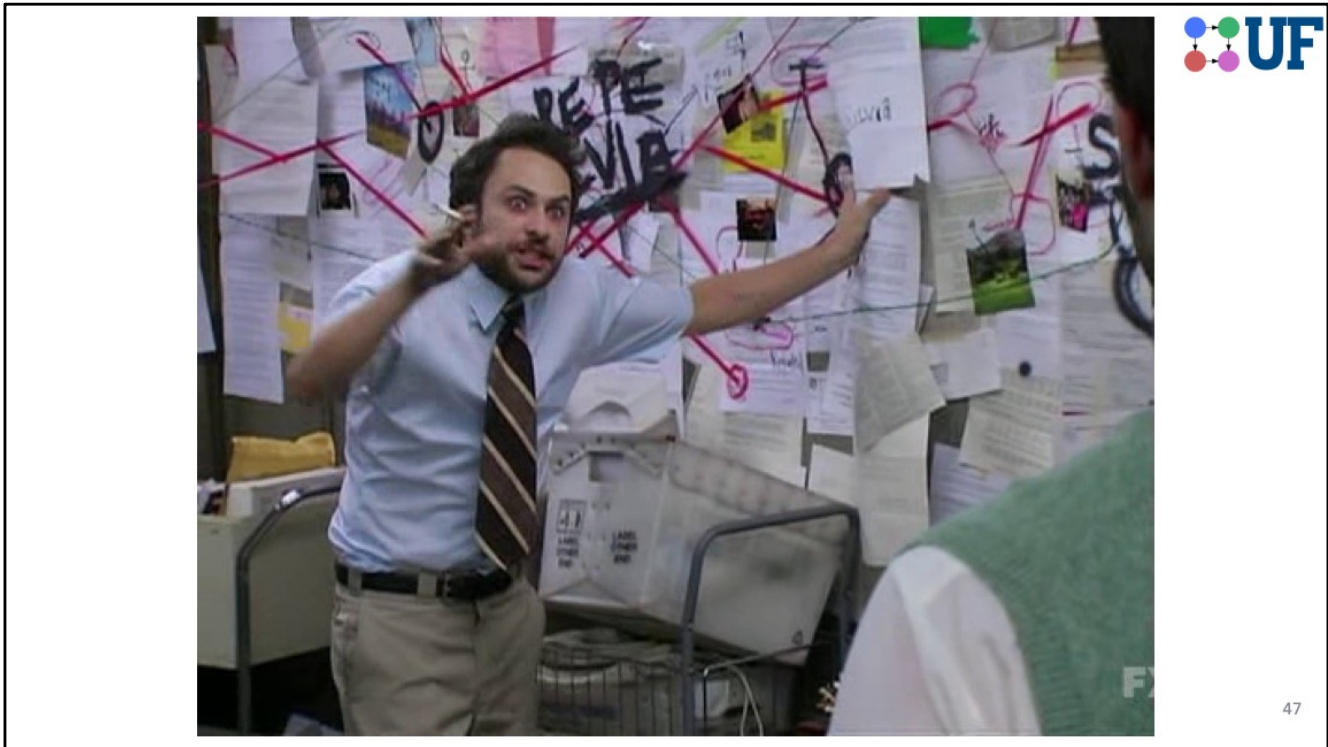
But I'll switch focus to some specific projects



I want to show three categories of projects with examples in each. I think of each of these projects as intersecting with each of the four topics I overviewed previously.



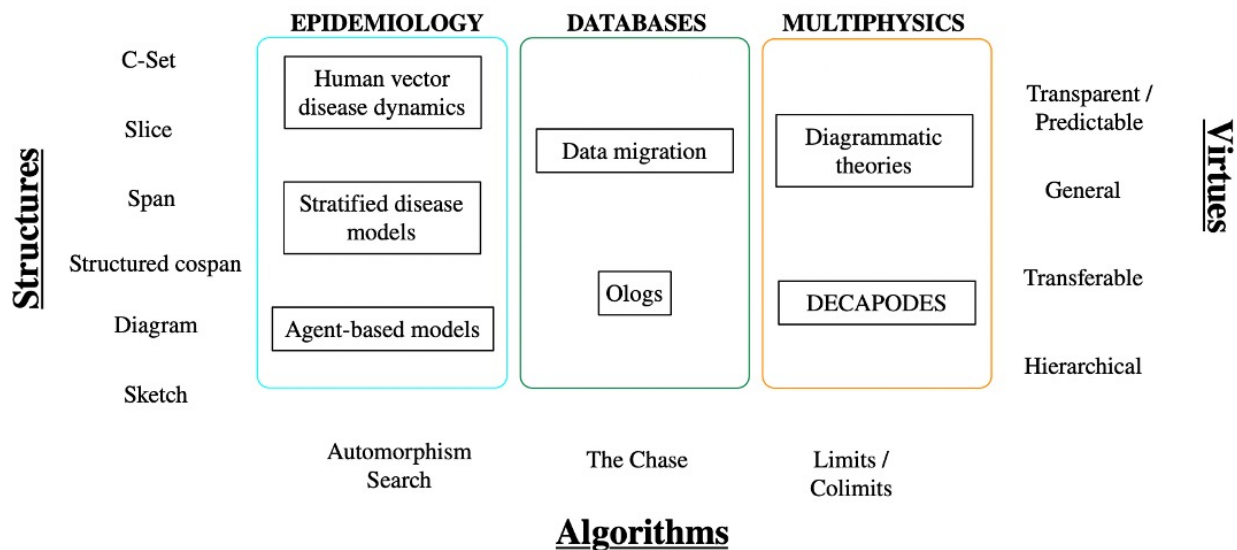
Looking at all the projects I want to talk about, we could draw all the connections.



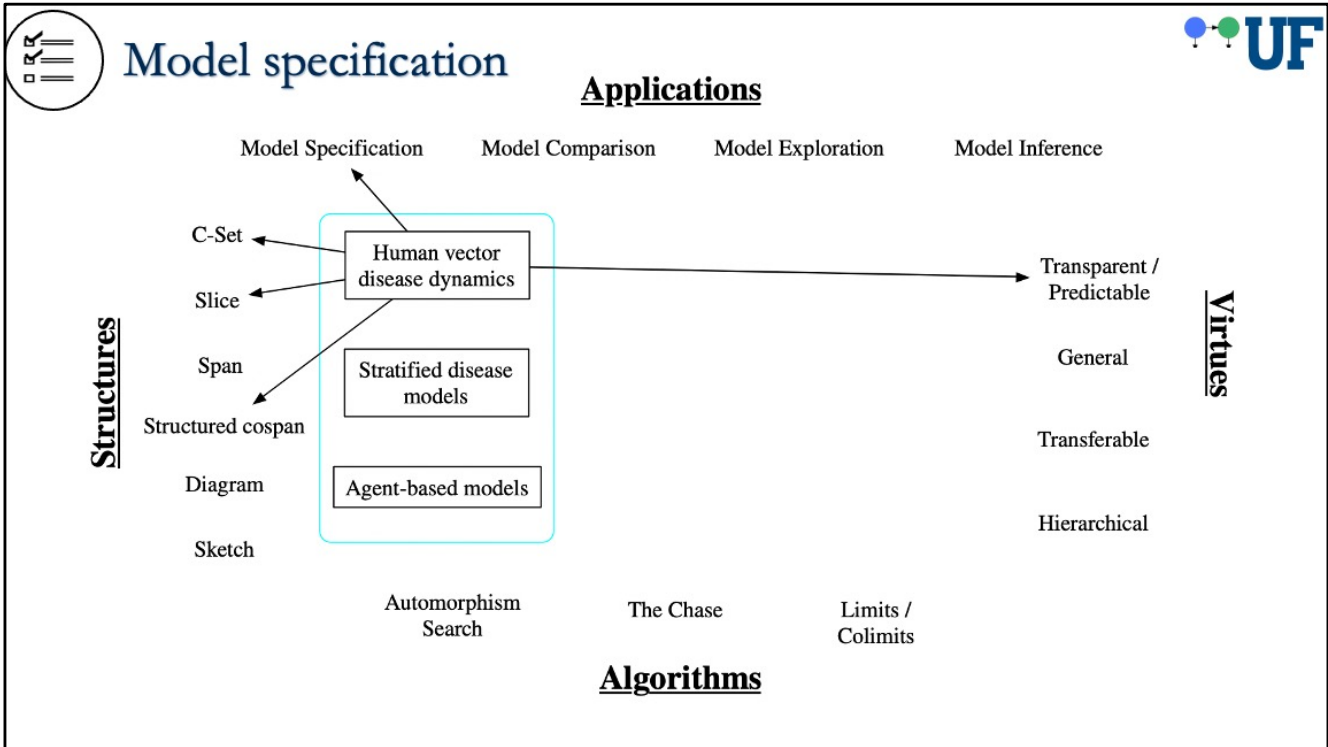
I've been strongly advised to not show all the connections at once.

Applications

Model Specification Model Comparison Model Exploration Model Inference



In fact, we will just focus on epidemiology.



The first thing I will talk about is modeling disease dynamics.

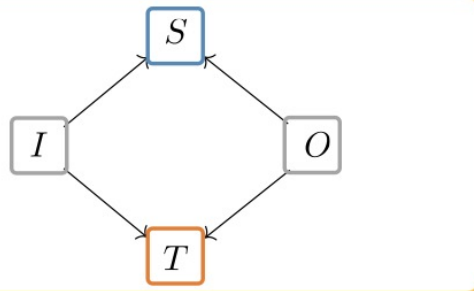


Specification: Structure and semantics

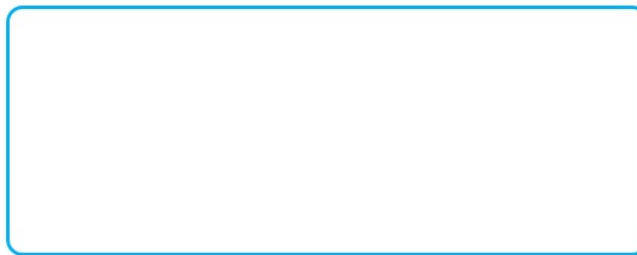


Petri

Schema:



Instance:



50

Now I'll reintroduce the concept of a Petri Net as a C-Set. A petri net is a nice way to separate the syntax (a kind of graph showing what reacts with what) from the semantics (which are dynamical systems).

We use a C-set with this shape of sets and functions to encode the fact that there are possibly many species flowing into and out of transitions.

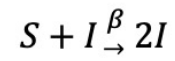
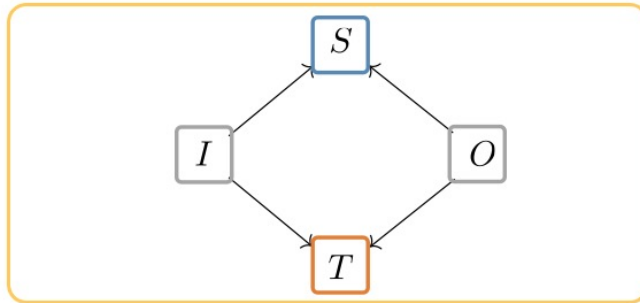


Specification: Structure and semantics



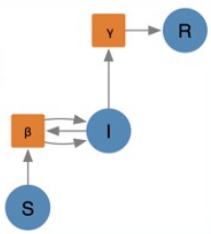
Petri

Schema:



Instance:

Input	I_S	I_T	Out	O_T	O_S
1	S	β	1	β	I
2	I	β	2	β	I
3	I	γ	3	γ	R



51

Here's SIR as an example on the bottom, with the equations that are represented by this model

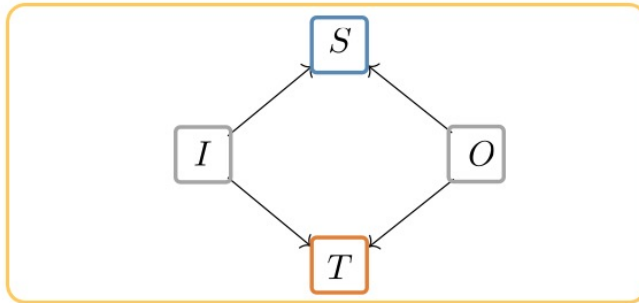


Specification: Structure and semantics

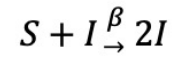


Petri

Schema:



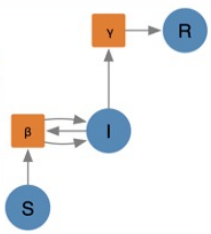
Petri Nets are C-Sets, can be given the semantics of mass-action kinetics



$$r_{\beta} = k_{\beta} * [S] * [I]$$

Instance:

Input	I_S	I_T	Out	O_T	O_S
1	S	β	1	β	I
2	I	β	2	β	I
3	I	γ	3	γ	R



$$r_{\gamma} = k_{\gamma} * [I]$$

We can assume mass-action kinetics, where reaction rates are given by a constant times the product of reactant concentrations.

So we could stop here and say we're done defining our class of models, but this is pretty unrestricted. There's not much scientific knowledge encoded here other than mass action kinetics.

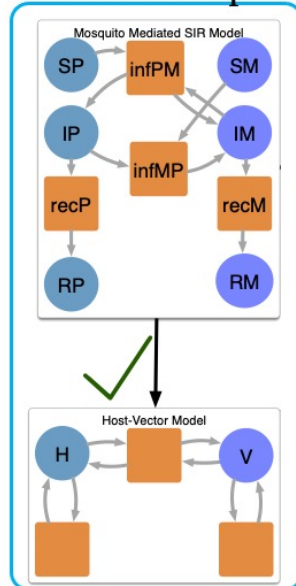
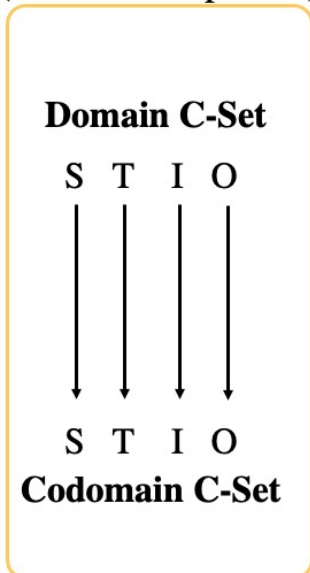


Specification: Combinatorial constraints



C-Set Transformation
(or 'homomorphism')

Example
Petri homomorphism



- C-Set morphisms are combinatorial data structures.
- They map elements of one C-Set onto elements of another.
- They must satisfy some laws (i.e. 'preserve the structure' of the domain)

← Codomain =
Constraints on
domain

Libkind et al. "An algebraic framework for structured epidemic modeling" (2022)

53

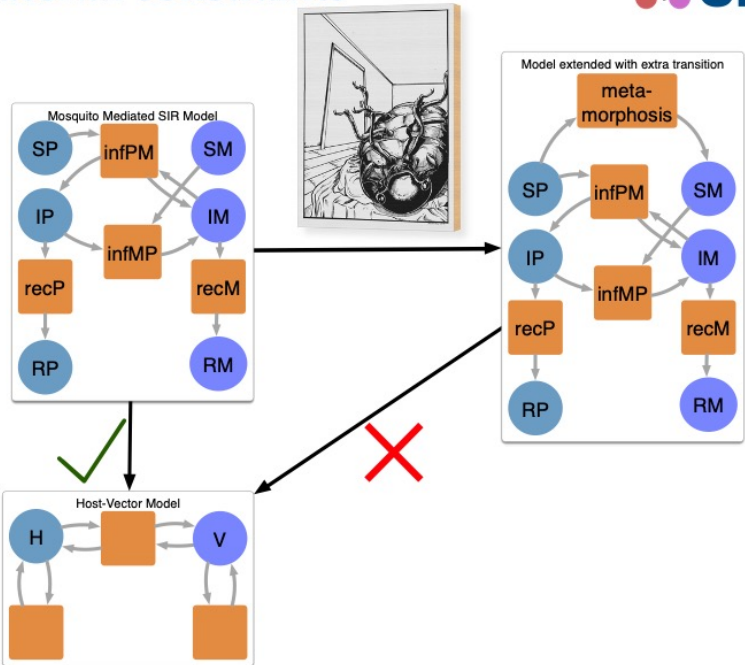
Now I want to return to those constraints I talked about earlier. Just like a C-set is a combinatorial data structure, the transformations between C-sets are also combinatorial data structures. A C-Set transformation is like a meta-level function which maps C-sets (which are networks of sets and functions) in a way that respects structure.

Suppose we have a particular petri net which encodes very high level information about the kind of petri nets we're interested in. This one is an example taken from a recent paper of Sophie / Evan / James.

By restricting our petri nets to those which have homomorphisms into this one, we're in effect demanding that anyone who provides us a Petri also must declare that their species are kinds of either Humans or Vectors, and that all transitions are either two humans interacting, two vectors interacting, or a human and a vector interacting.



Specification: Combinatorial constraints



Demanding that a C-Set morphism exists into some fixed model is a constraint.

Libkind et al. "An algebraic framework for structured epidemic modeling" (2022)

54

Why do I call this a constraint? Well not every petri net has a homomorphism into it. The constraint can be efficiently checked, so it's not possible for someone to submit a new model that has a transition converting a human into a vector. We rule out that (and many more nonsensical Petri nets) merely by changing our model type from Petri net to this homomorphism type (which is called a slice in category theory lingo).

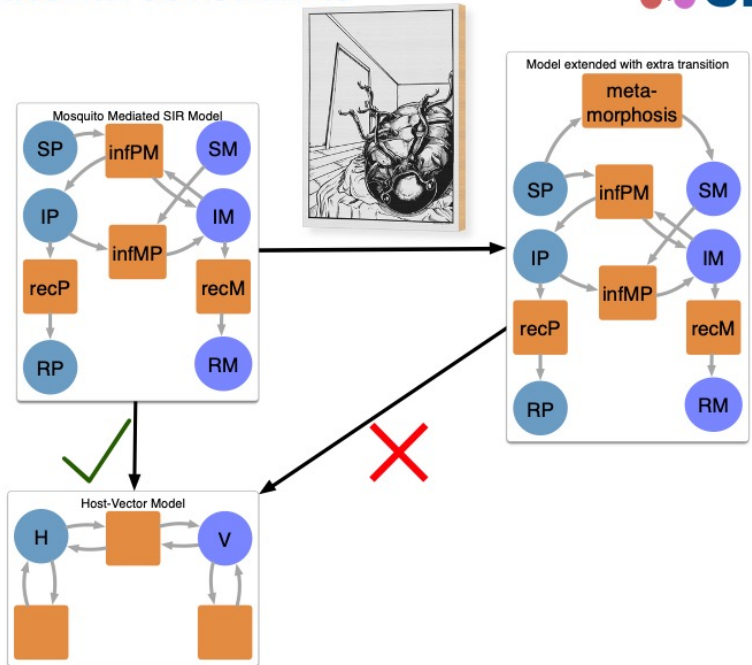


Specification: Combinatorial constraints



- Homomorphisms effectively assign “types” to elements of the domain model.

Demanding that a C-Set morphism exists into some fixed model is a constraint.



Libkind et al. "An algebraic framework for structured epidemic modeling" (2022)

55

Not only have we constrained our models, but we obtain a notion of typed elements that is completely transparent in the model: we know which states are humans and which are vectors.

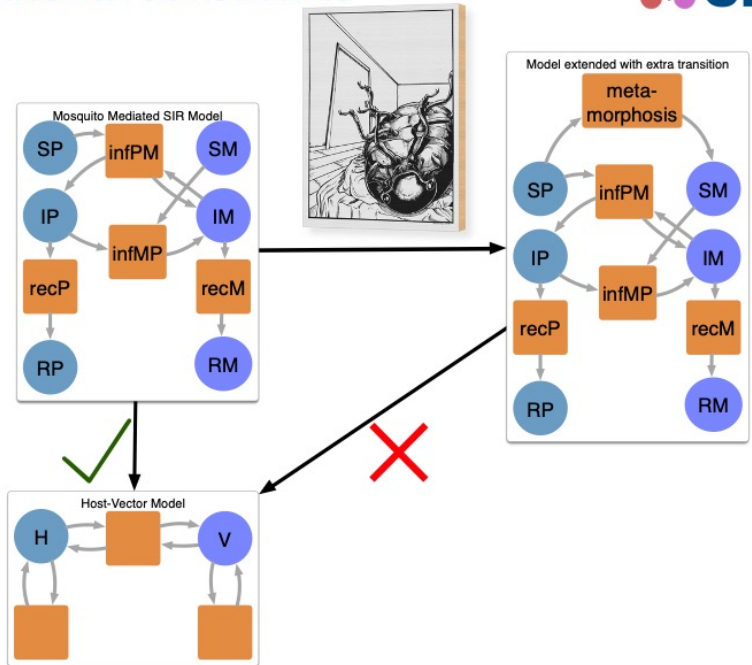


Specification: Combinatorial constraints



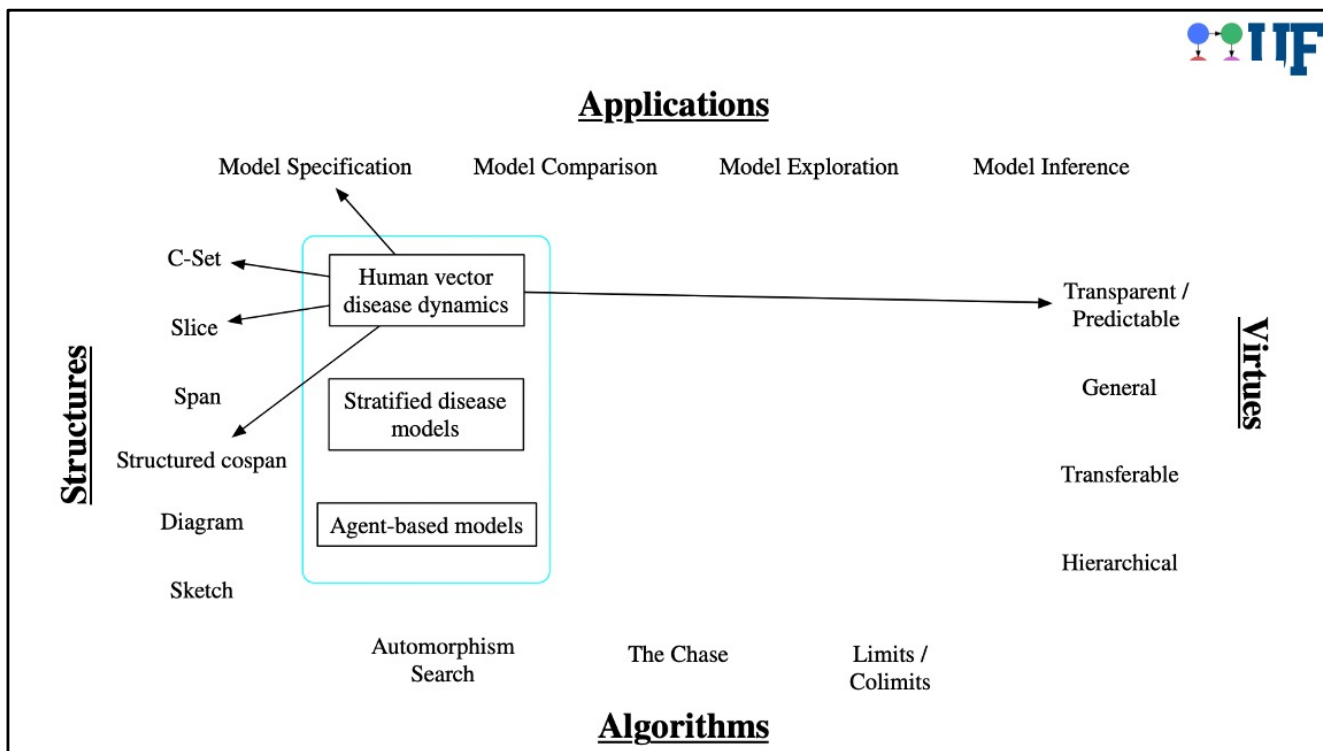
- Homomorphisms effectively assign “types” to elements of the domain model.
- Homomorphism search is type inference

Demanding that a C-Set morphism exists into some fixed model is a constraint.



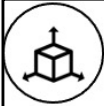
Libkind et al. "An algebraic framework for structured epidemic modeling" (2022)

And one useful feature of having homomorphism search implemented generically (for all C-sets) is that we get for free an algorithm which takes a partially labeled Petri net and tries to find whether zero, one, or many type type assignments. This is good because it can be tedious to specify all this data manually.



So in summary, we enable the representation of models encoding precise domain knowledge by using the combinatorial structure of slices.

I actually could use this story to demonstrate all four of the virtues here, but the key one to stress is the fact that model has become transparent and predictable from representing it combinatorially.



Model exploration

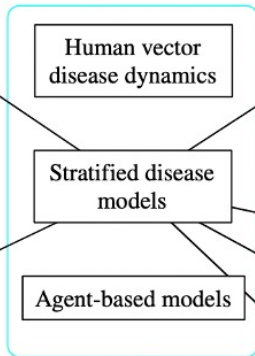


Applications

Model Specification Model Comparison Model Exploration Model Inference

Structures

C-Set
Slice
Span
Structured cospan
Diagram
Sketch



Automorphism Search

The Chase

Limits / Colimits

Algorithms

Virtues

Transparent / Predictable
General
Transferable
Hierarchical

Now I want to return to the model exploration problem.

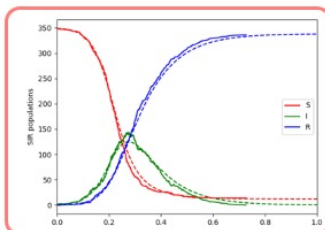
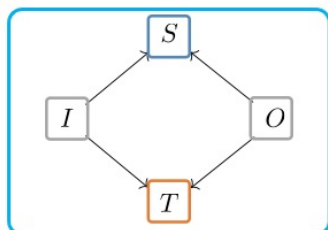
The key structure we'll be invoking is a diagram. In particular, we'll look at diagrams in the category of Petri Nets. (22:30)



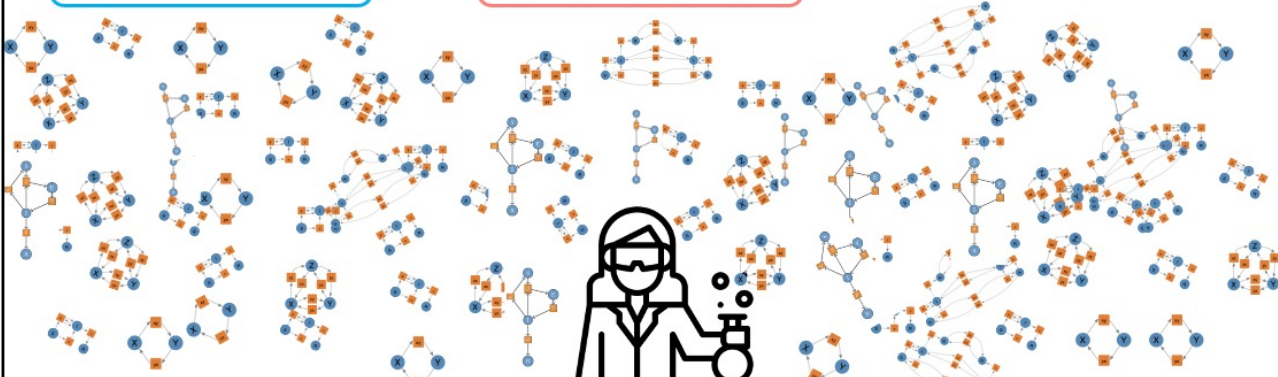
Exploration: Problem statement



“Given *this* class of models and *this* data ...what model best represents the data?”



We *must* impose some structure (implicit or explicit) on the space of models in order to navigate.



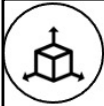
59

But first, what is our goal? Whatever problem we're trying to solve with the model, a scientist is confronted with the fact that there are too many possible models to do any kind of naive search.

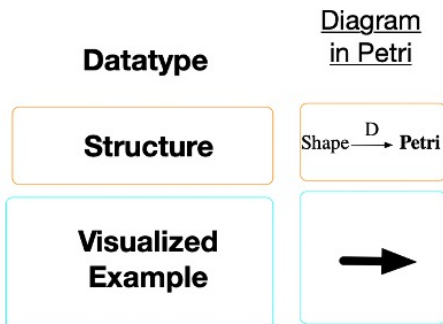
And it's hard to explore if we don't have an actual structure through which allows us to visualize a space of Petri nets (rather than this big grab bag soup that comes from considering Petri nets as living in a big Set rather than something with more structure).

So I'll now introduce the category theoretic notion of diagrams as a potential solution to creating a space of models rather than a bag of models.

Diagrams give us a tool for cutting through this space in a meaningful (not brute force) way.



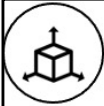
Exploration: diagrams as explicit model spaces



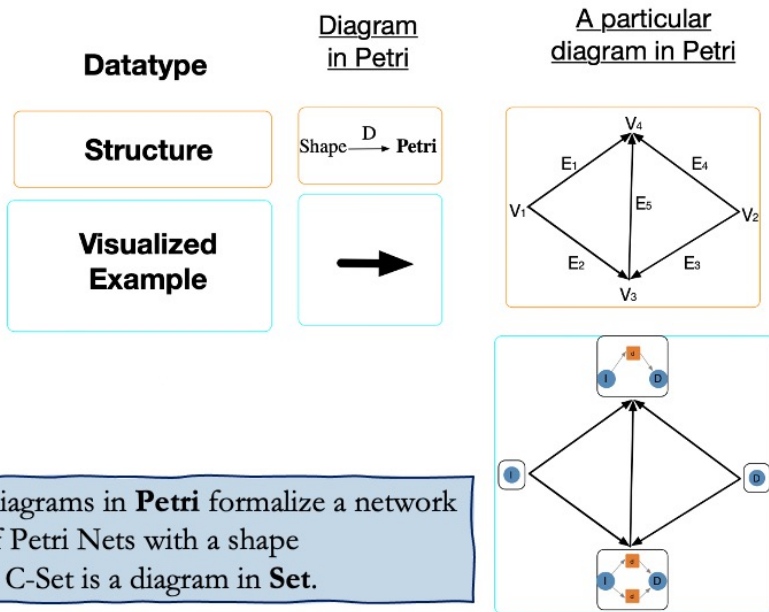
- Diagrams in **Petri** formalize a network of Petri Nets with a shape
- A C-Set is a diagram in **Set**.

A diagram in Petri is a combinatorial data structure built on top of C-Sets. In essence, we have a shape which is a schema just like “C” was a schema for C-sets (it’s like a directed graph with equations). And instead of the values we put on that structure being Sets and Functions, we put Petri Nets and Petri net morphisms.

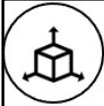
Note, given this definition, a C-set is itself tantamount to a diagram in Set.



Exploration: diagrams as explicit model spaces



Here's an example defining a triangle-pair-shaped network of Petri nets.



Exploration: diagrams as explicit model spaces



Datatype

Structure

Visualized Example

Diagram
in Petri

Shape \xrightarrow{D} Petri



A particular
diagram in Petri

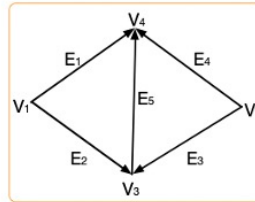
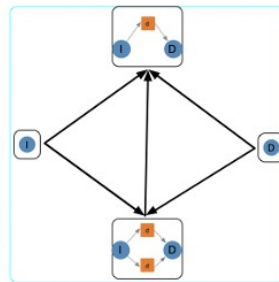
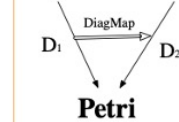


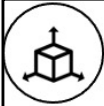
Diagram
Morphism
in Petri

Shape₁ $\xrightarrow{\text{ShapeMap}}$ Shape₂



- Diagrams in **Petri** formalize a network of Petri Nets with a shape
- A C-Set is a diagram in **Set**.

That's a very obvious generalization to make, but the subtlety lies in what the right choice of morphisms are between these objects. For C-sets, we only ever look at morphisms where C is fixed, but now this category of diagrams can have changing structure and changing data at the same time!



Exploration: diagrams as explicit model spaces



Datatype

Structure

Visualized Example

Diagram
in **Petri**

Shape \xrightarrow{D} **Petri**



A particular
diagram in Petri

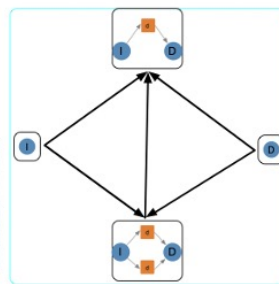
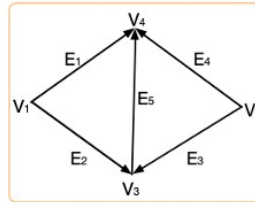
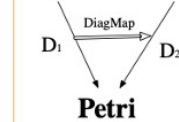


Diagram
Morphism
in **Petri**

Shape₁ $\xrightarrow{\text{ShapeMap}}$ Shape₂



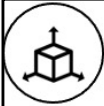
Petri

Don't worry!

- Diagrams in **Petri** formalize a network of Petri Nets with a shape
- A C-Set is a diagram in **Set**.

It's a bit too much data to visualize, but trust me that this is the right notion of morphism for our purposes.

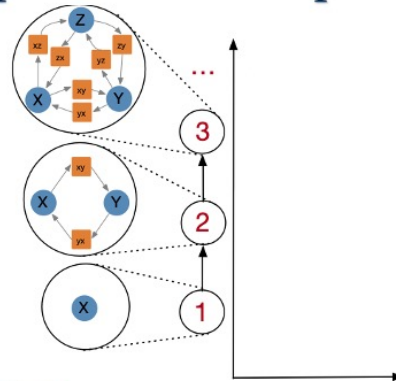
But now I want to address the question of why it was actually useful to recast our vague notion of a "model space" into something precise like this - what did we gain? Well, category theory has a lot to say about things you can do with diagrams, and it turns out some of those are useful at recasting our other vague notions into concrete algorithms that can be implemented.



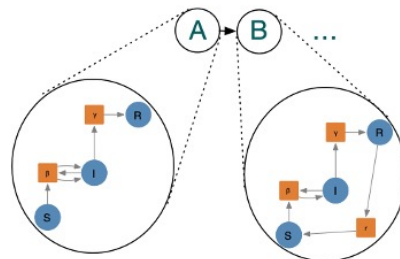
Exploration: product model space



“City dimension”



Products of diagrams take different ‘dimensions’ and yield a product model for each combination of elements from each dimension.



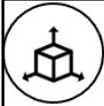
“Disease dimension”

64

One thing we get is a way of combining small model spaces into larger ones. This is a reasonable problem to have, where the scientist has a few ideas of “dimensions” so to speak that they want to explore along, and they want their models to take features from each of these dimensions.

Here’s an example of a pullback in the category of diagrams (though because pullbacks are a kind of generalization of products, I also call it a product of model spaces).

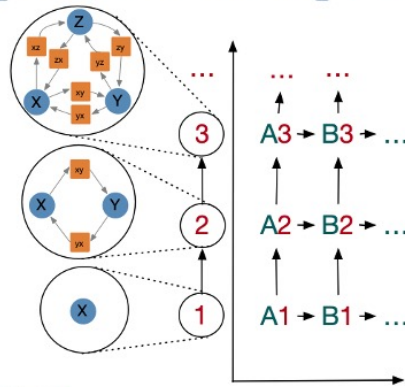
On the vertical axis, we have a 1-city model, a 2-city model, and so on. And on the horizontal axis we have some sequence of disease dynamics.



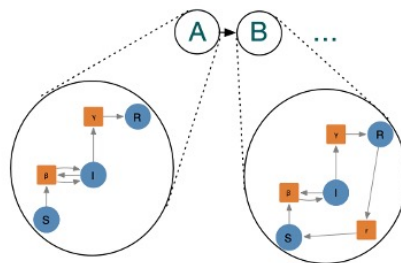
Exploration: product model space



“City dimension”



Products of diagrams take different ‘dimensions’ and yield a product model for each combination of elements from each dimension.



“Disease dimension”

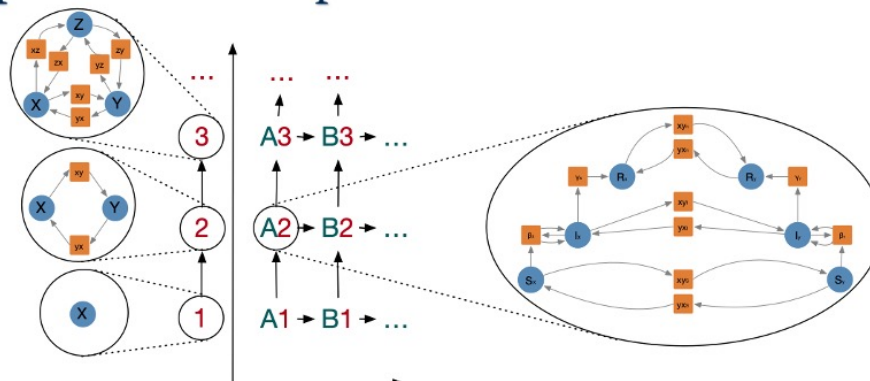
The resulting model space would look like this at the shape level. We’ll pick one of these models to see what element of Petri that vertex got mapped to.



Exploration: product model space



“City dimension”



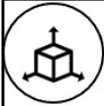
Products of diagrams take different ‘dimensions’ and yield a product model for each combination of elements from each dimension.

“Disease dimension”

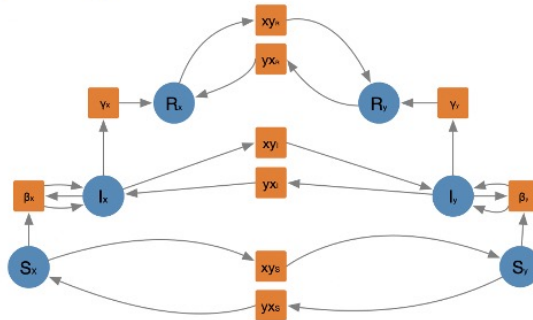
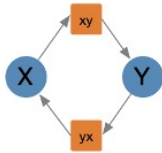
66

You see here we have the same disease dynamics occurring in two cities that interact with each other.

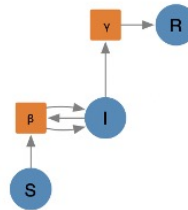
If this process were coded up in a script, it would look like a nested for-loop for each dimension. That works for a one off product, but it doesn’t compose well with other ways of combining model spaces and the resulting product models do not have provenance.



Exploration: model space provenance

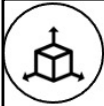


Products of diagrams leave a record of where elements of a product Petri Net came from, which can be used for parameter estimation.

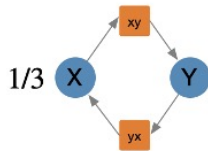


67

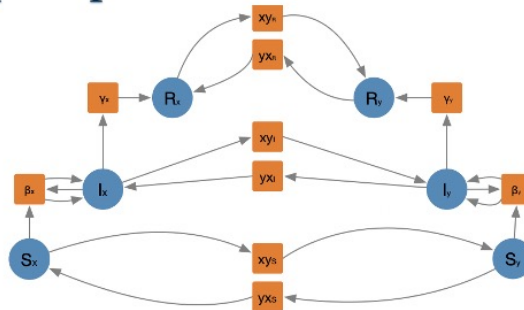
Computation of pullbacks gives a complete account of where features in a composite model come from. This can be useful in many ways. The one I will highlight is making initial guesses for parameters.



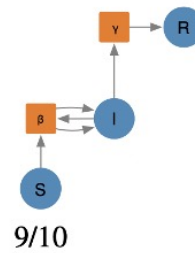
Exploration: model space provenance



1/3



Products of diagrams leave a record of where elements of a product Petri Net came from, which can be used for parameter estimation.



9/10

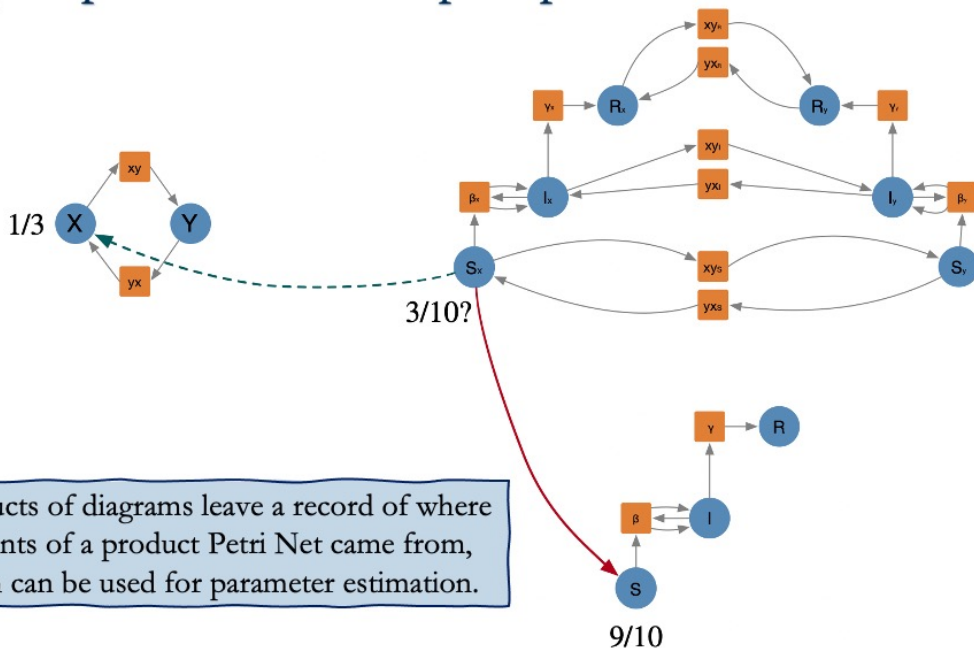
68

Suppose we fit SIR to our data and determine initially 9/10 people are Susceptible, and we collect some demographic data that 1/3 of people are in City X.

We really ought to deduce that, initially, 3/10 people are susceptible in City X in the product model.



Exploration: model space provenance



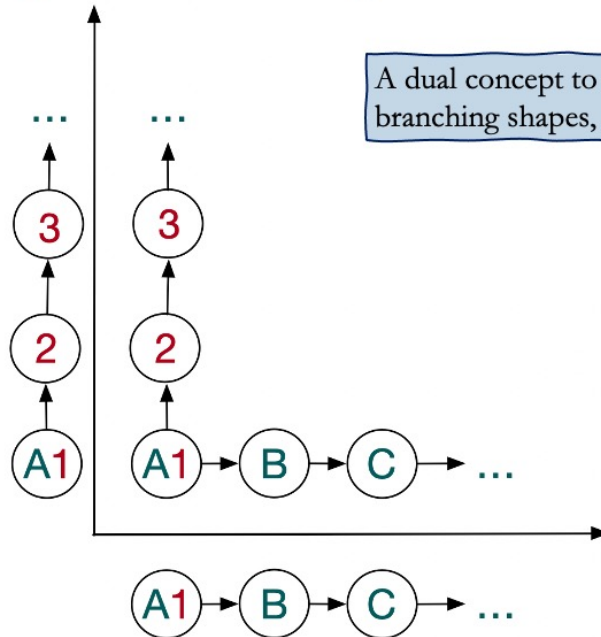
69

And you'll see that the data of a pullback actually tells us where that state was derived from, so we can compute this value!

These types of computations assume a kind of statistical independence of the dimensions, but it is still useful as an initial guess for iterative algorithms.

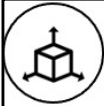


Exploration: pushout model space

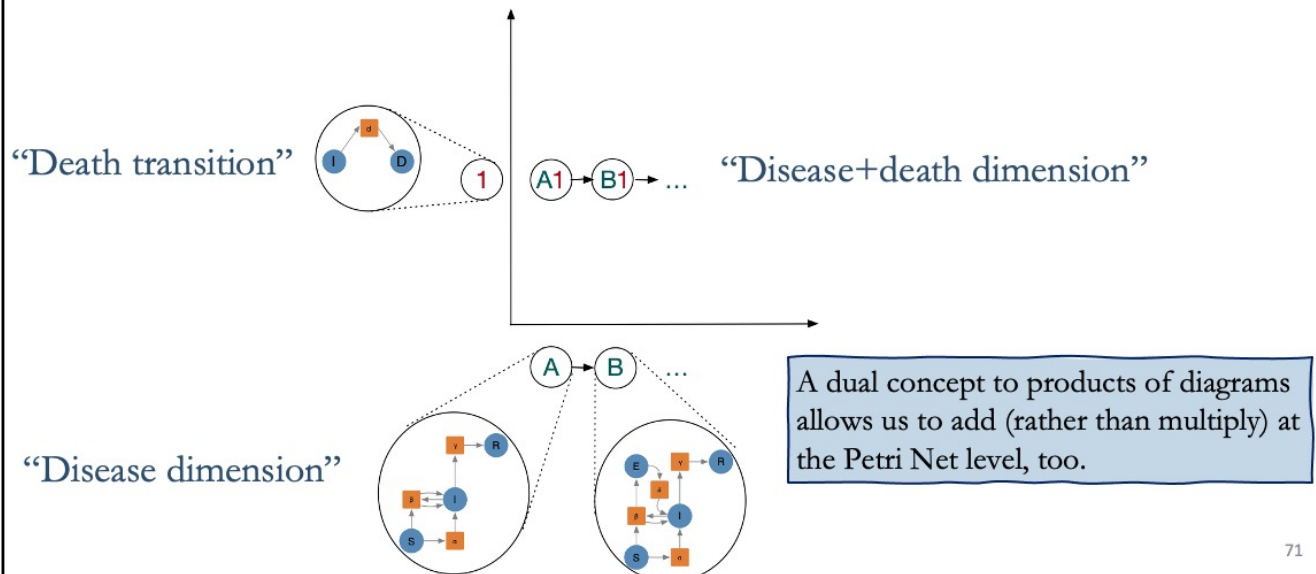


70

Just like we can take a product of two model spaces, we can also take a sum, which would produce more of a branching shape than a grid shape here.

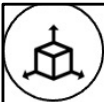


Exploration: pushout model space

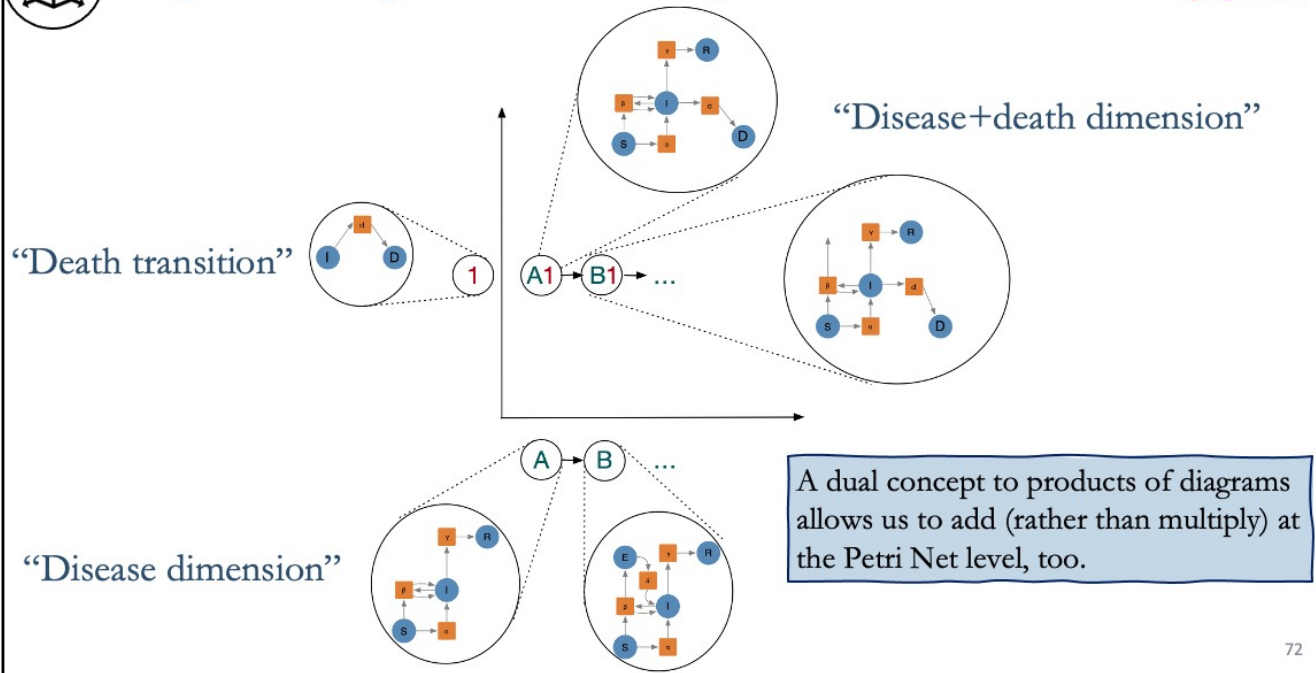


71

And this other example which I'll gloss over is an example of gluing together things at the Petri-net level, not the shape level. So we have an Infected to Death transition that we attach to a sequence of models in the bottom model space.



Exploration: pushout model space

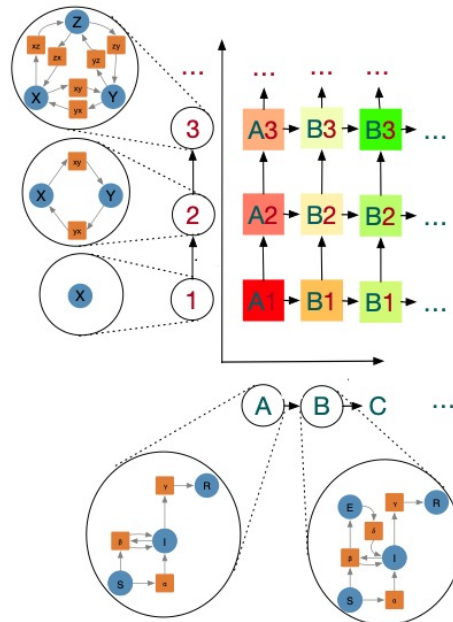


72

And specializing this very basic CT idea of a pushout in the context of diagrams computes the correct thing for free.



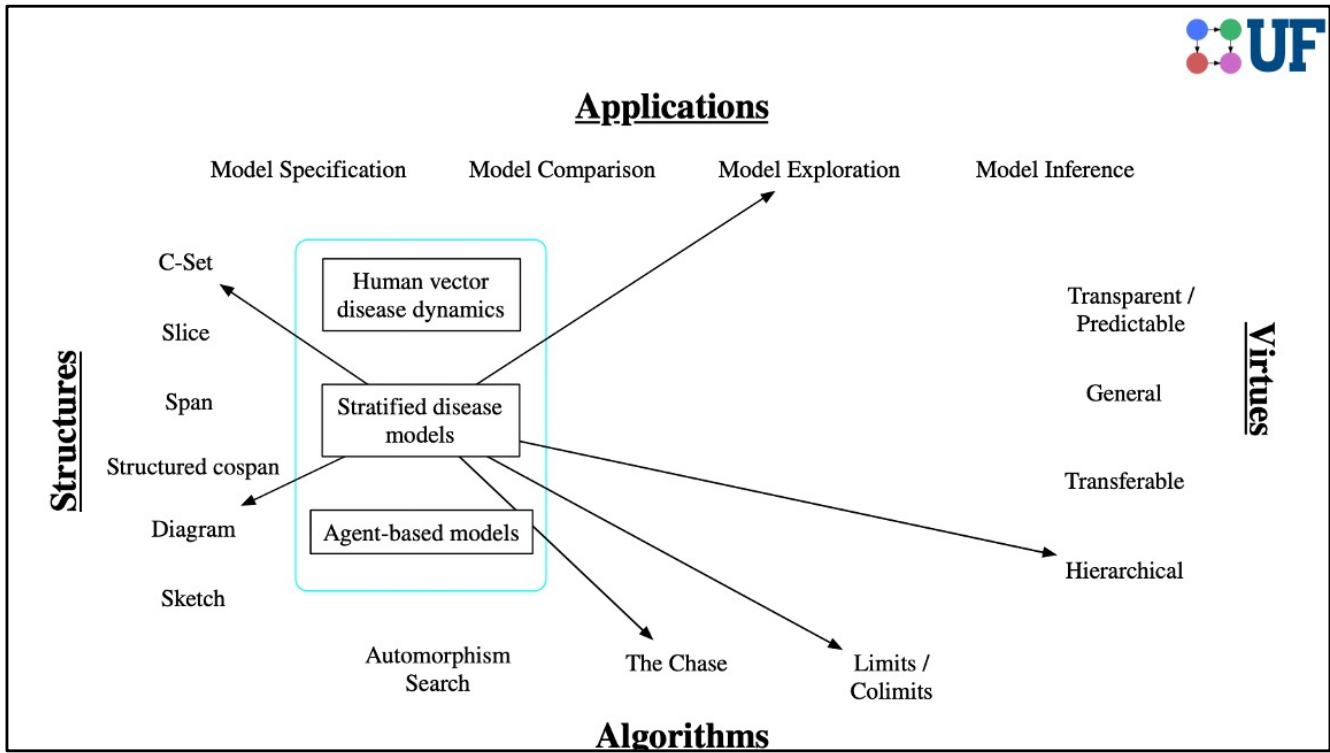
Exploration: model selection



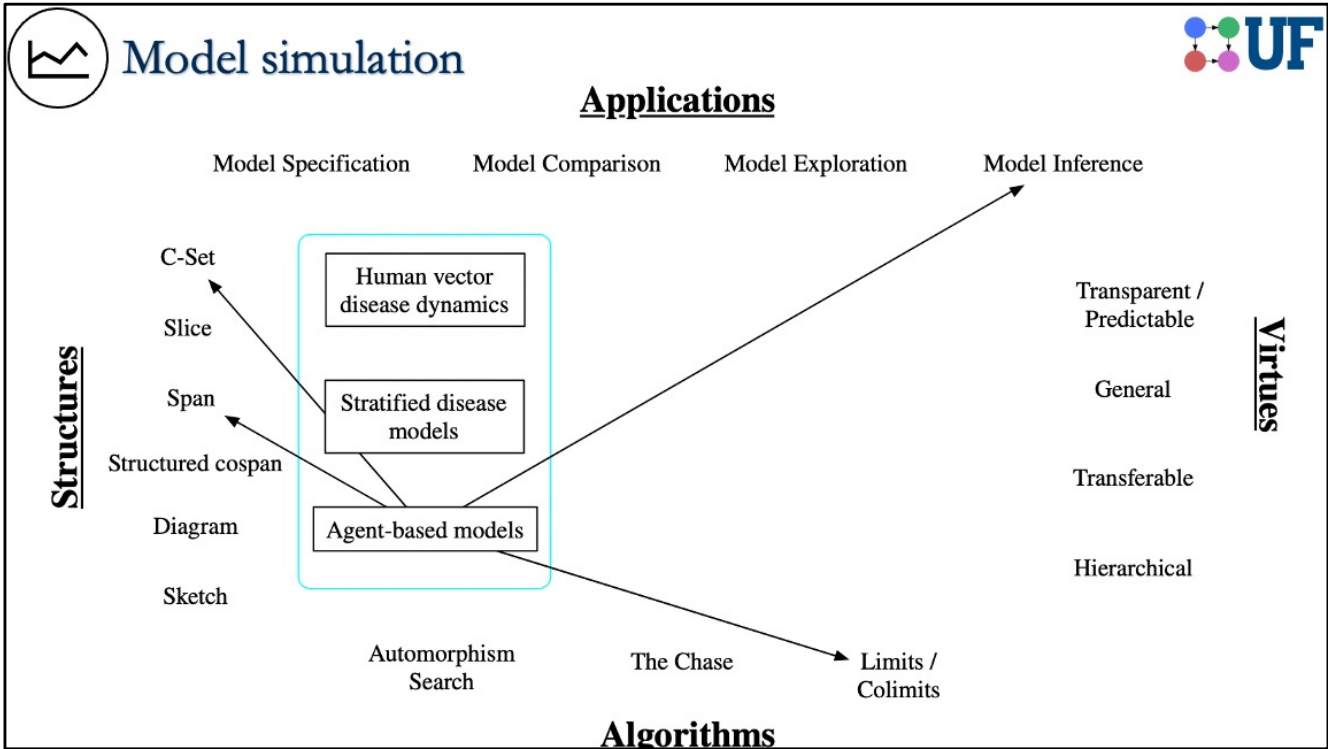
After defining a model space, one can efficiently traverse it to find optimal models.

73

Ok, that was two very simple examples of two very basic ways to combine model spaces into larger ones. I want to switch focus to the problem of having defined the composite model space, you are now trying to select to pick which model is your favorite. Say you have a function that can give each model a score. These basic compositions introduce a geometry on the model space; we can also get a notion of gradient. We can navigate this space efficiently, rather than by brute force.



So that was the second project I wanted to highlight.



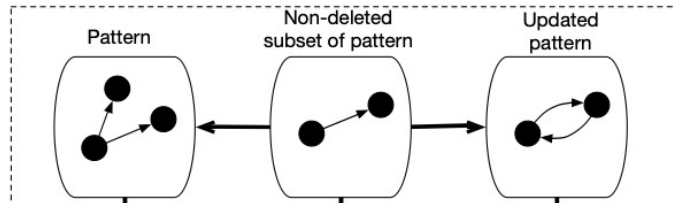
Lastly I want to talk about a model inference problem. Often the semantics of a scientific model are some kind of dynamical system. We'll focus on a flavor of discrete dynamical systems called rewriting systems.



Simulation: rewrite rules



- In general, there is a **pattern match**, a **deletion**, and an **insertion**.
- Categorical rewriting captures all three with a *span* of morphisms



A recipe for performing a *dynamic* update can be represented by a *static* data structure. The same few lines of code performs rewriting on *any* C-Set.

76

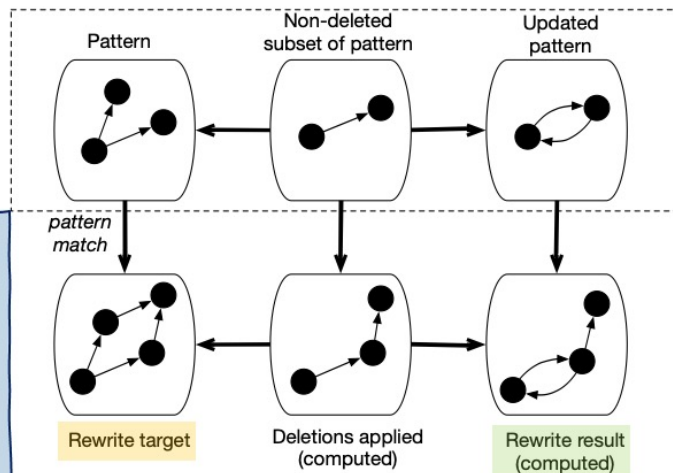
Main mechanism of discrete update in a simulation is applying a rewrite rule. The structure of a rule is to have a pattern match, a deletion, and an insertion. We break the rule down into a pattern that gets matched onto our model that we want to rewrite. We have a portion of the match that is preserved and we have a third structure which is what gets substituted into the location that was matched.



Simulation: rewrite rules



- In general, there is a **pattern match**, a **deletion**, and an **insertion**.
- Categorical rewriting captures all three with a *span* of morphisms

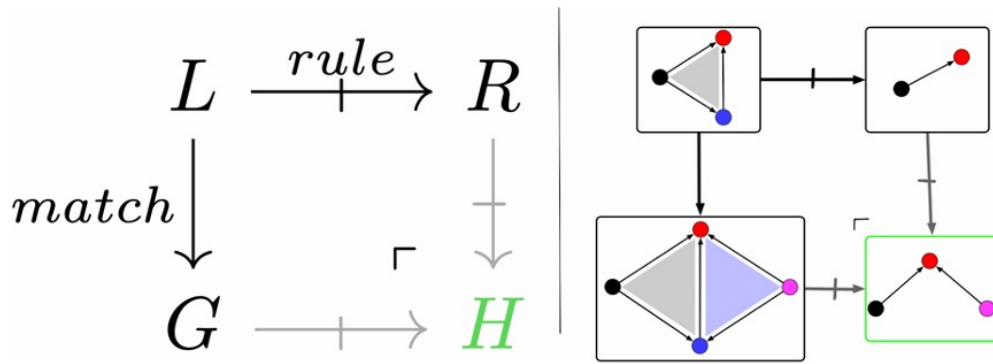


A recipe for performing a *dynamic* update can be represented by a *static* data structure. The same few lines of code performs rewriting on *any* C-Set.

77

So you see here how this process applies to the bottom left graph and yields as the final result the graph on the bottom right. The actual implementation is defined in terms of categorical operations that make sense for all ACsets.

This is a lot to throw at you very quickly, but the conceptual point here is that we were able to express this inherently DYNAMIC process of rewriting with a purely STATIC structure, which is those three structures there in the box. We did not need to write a function in code or a differential equation.



SPO is a different semantics that can be given to a rewrite rule. It enables "cascade delete".

What I just showed was double pushout rewriting (DPO). Single Pushout Rewriting is different from DPO because when you delete something, you implicitly delete everything that refers to it. This is like CASCADE DELETE in databases.

I shown an example of this using a different C-Set from Petri Nets just because it is easier to visualize. We delete the blue vertex, but the lavender triangle is referring to the blue vertex so it gets implicitly deleted, even though it's far removed from the part of the model that got matched by the pattern.

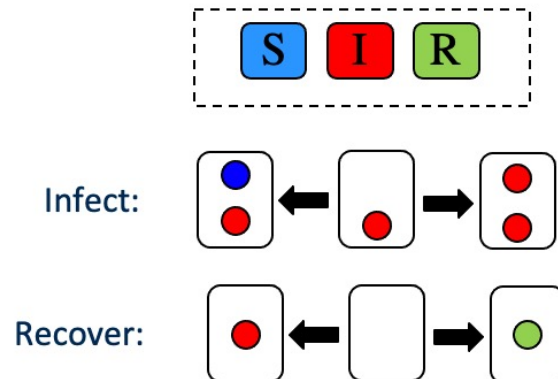
Here is another example of the same syntax (the rewrite rule) being given multiple semantics.



Simulation: SIR agent based model



We can give Petri Nets rewrite rule semantics rather than ODE semantics.



79

So taking the world state to be represented by a C-Set with just three objects and no morphisms, we can recapture the dynamics of SIR by these two rewrite rules.

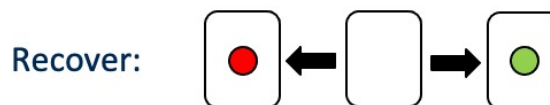
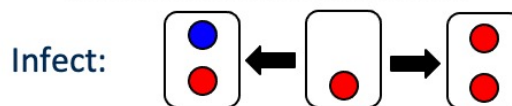
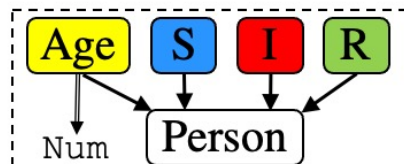
These could be generated automatically from the Petri Net representation.



Simulation: Stratified agent based model



We can give Petri Nets rewrite rule semantics rather than ODE semantics.

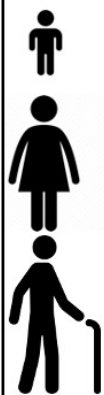


80

Moving beyond that toy example, we could add a numerical age to the schema.

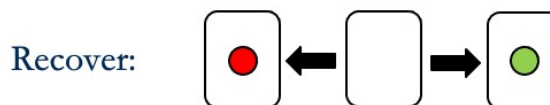
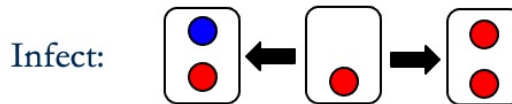
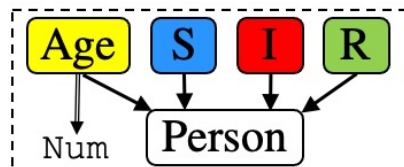


Simulation: Stratified agent based model



C_{11}	C_{12}	C_{13}
C_{21}	C_{22}	C_{23}
C_{31}	C_{32}	C_{33}

We can give Petri Nets rewrite rule semantics rather than ODE semantics.



81

We might have some function which takes two ages and computes the probability of that infection happening.

One type of function would be to break down demographics into bins, and n^2 numbers characterize each group's propensity to infect another.






Simulation: Stratified agent based model

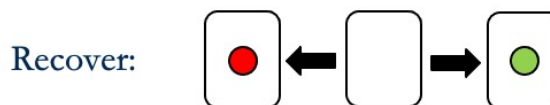
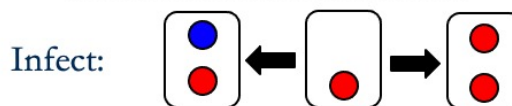
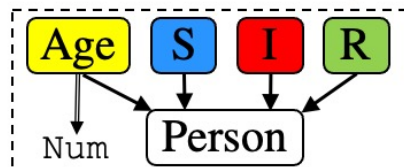


We can give Petri Nets rewrite rule semantics rather than ODE semantics.



- n^2 rewrite rules
- 1 rewrite rule + 1 line of 'arbitrary' code

	C_{11}	C_{12}	C_{13}
	C_{21}	C_{22}	C_{23}
	C_{31}	C_{32}	C_{33}



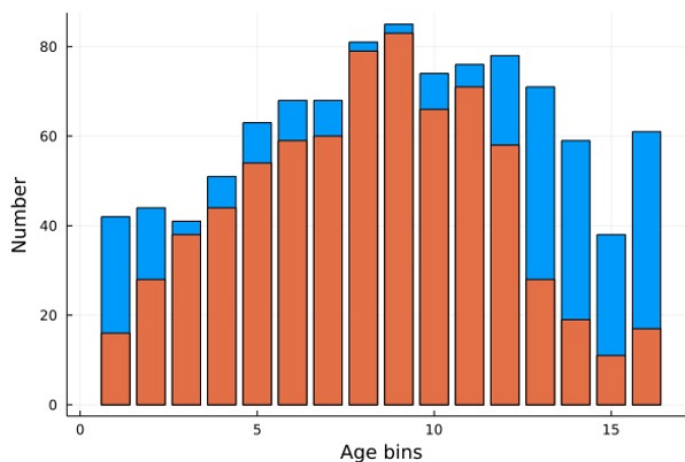
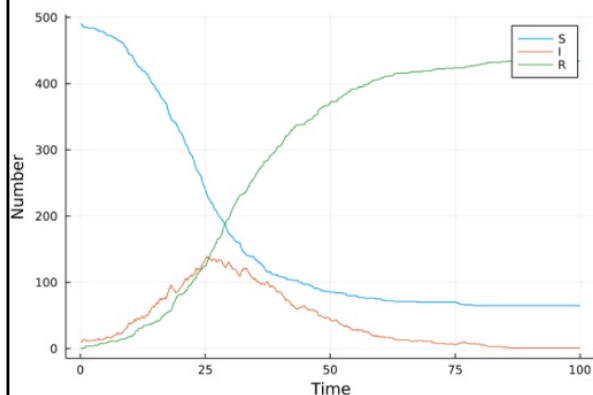
82

Here there is an interesting tradeoff where having n^2 rewrite rules is completely transparent to a computer, but one rewrite rule + 1 line of arbitrary code feels more human interpretable.

But we're in a position to pick which design makes the most sense.



Simulation: Age-stratified Results



We implemented a prototype agent based model framework for any C-set model.

Wu et al. "Individual-based models based on graph rewriting" (2022).

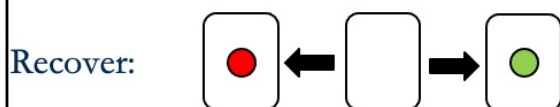
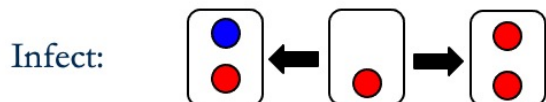
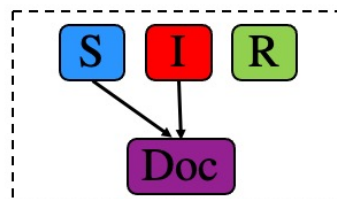
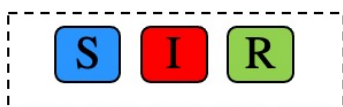
So we used a real contact matrix based on research in Taiwan, and Sean Wu at the institute of Health Metrics and Evaluation was able to reproduce earlier results from Individual.jl by using C-set rewriting as the core engine.

In the plot on the right, the orange fraction of the bar of an age group shows what fraction got infected over the course of the simulated pandemic.

You can see the elderly and young are relatively insulated from the working age population which bears the brunt of infections.



Simulation: Algorithm transferrability



- Rewrite rules express algorithm as data.
- This means we can migrate our entire algorithm as our assumptions change.

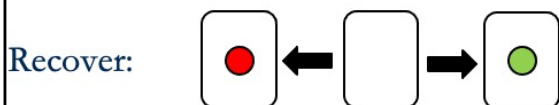
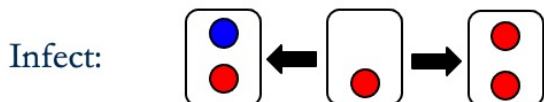
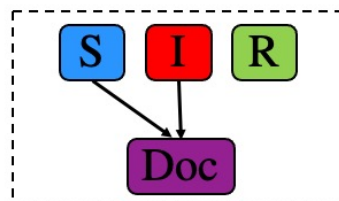
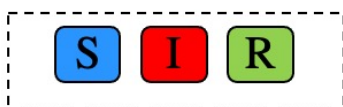
84

Here's one reason why one might prefer the machine interpretable route. I'm going to demonstrate the transferrability virtue, which we obtain by having our algorithm itself expressed combinatorially.

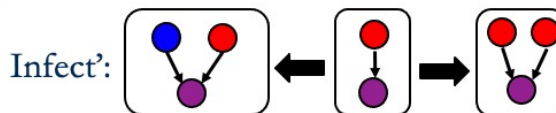
Suppose we decide to update our schema like so, where we model the fact that susceptible and infected people each are assigned a doctor from some set of doctors.



Simulation: Algorithm transferrability



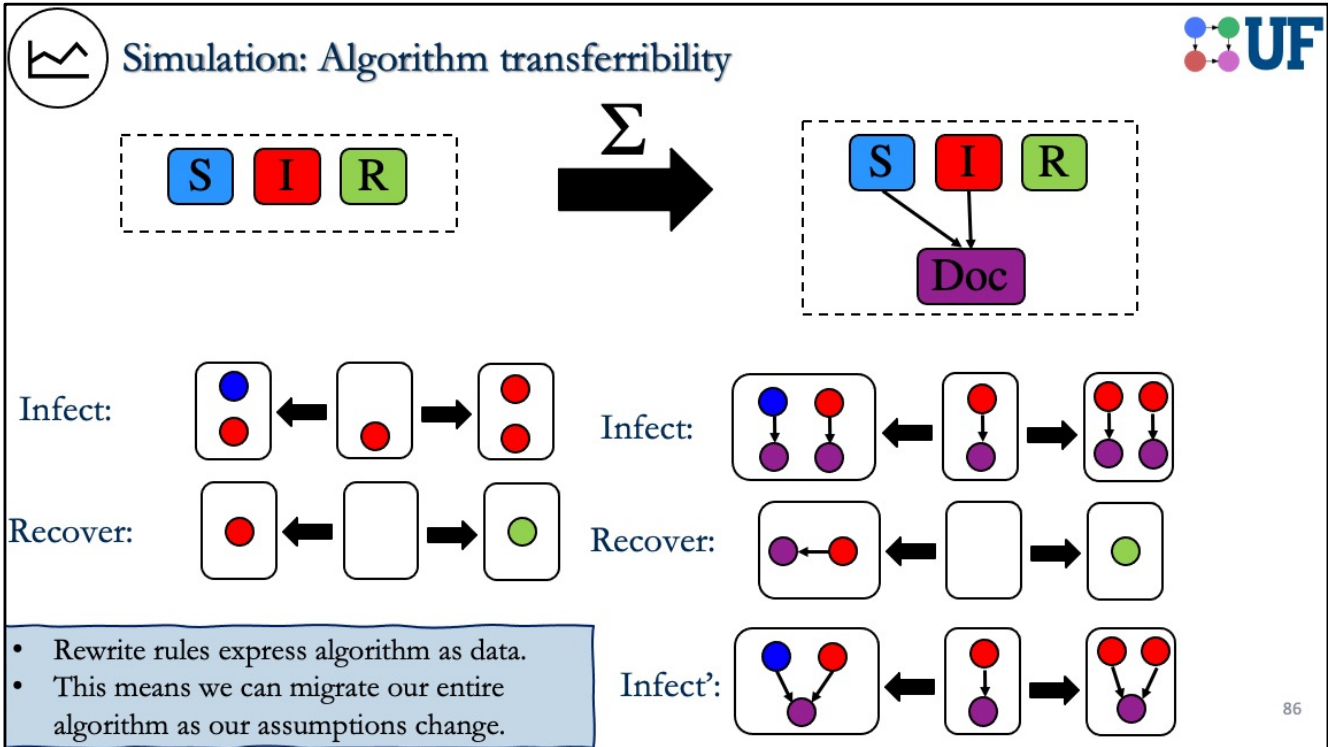
- Rewrite rules express algorithm as data.
- This means we can migrate our entire algorithm as our assumptions change.



85

This might be because we deliberately want to introduce an alternate mode of infection that only applies to people who share the same doctor (maybe in the waiting room, or the doctor is a vector), which could be written like this.

But we might have a lot of infrastructure built around our old SIR model. How do we "import" this into our new schema?



It turns out there is a very obvious functor from the left category into the right one, and by pushing the old rewrite rule into the new schema, we obtain the correct interpretation of it, which freely assigns a doctor to each person in the rewrite rule.

This is an example of migrating an algorithm, and likewise when we express queries combinatorially (e.g. by representing queries as homomorphisms), we can migrate our analysis as well.

This should be pretty mindblowing to people who think of algorithms as only expressible as general purpose code.

Takeaways



$\forall x, y: P(x) \rightarrow Q(y)$

$$\frac{\partial^2 u}{\partial x^2} + \frac{\partial^2 u}{\partial y^2} = 0$$

```
for step in range(1,50):
    print("Step ", step)
    rate = 0.5 * H2**2 * O2
    H2 -= 2*rate*dt
    O2 -= rate*dt
    H2O += rate*dt
    results.append(H2O)
```

- **Math, code, and formal logic** powerful enough to represent any kind of knowledge
- Easy to feel these as the only possible kinds that are formal.
- Strong expressive power trades off with ability for humans and computers to infer things

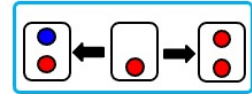
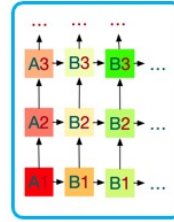
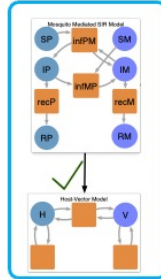
So despite the variety of examples, I think there was a central thesis to this talk. We're trying to present a different paradigm to scientists who think that math, code, or formal logic are the only ways to rigorously represent their domain knowledge. These are incredibly flexible and powerful tools, but that is a double edged sword because they are very difficult to reason about, for a human or especially a computer.

Takeaways

$\forall x, y: P(x) \rightarrow Q(y)$

$$\frac{\partial^2 u}{\partial x^2} + \frac{\partial^2 u}{\partial y^2} = 0$$

```
for step in range(1,50):
    print("Step ", step)
    rate = 0.5 * H2**2 * O2
    H2 -= 2*rate*dt
    O2 -= rate*dt
    H2O += rate*dt
    results.append(H2O)
```



- **Math, code, and formal logic** powerful enough to represent any kind of knowledge
- Easy to feel these as the only possible kinds that are formal.
- Strong expressive power trades off with ability for humans and computers to infer things

- Scientist's interface: intuitive + rigorous diagrams
- Model is **data**, not arbitrary math / code
- Better readability + inference for humans (and especially for machines)
- Model domain categorically (e.g. as a C-set) to gain access to powerful abstractions and algorithms

On the other hand, combinatorial data is inherently visualizable, making the barrier to rigorous modeling lower.

We showed lots of examples of why it's nice when computers can automatically reason about your knowledge, such as migrating information and algorithms from one domain to another in an automated way.

We're particularly interested in C-sets which have a low barrier to entry, wide generality, and many powerful abstractions and algorithms to use with them.

Thanks!

Evan Patterson



Sophie Libkind



David Spivak



T. Hanks



James Fairbanks



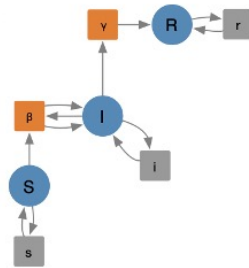
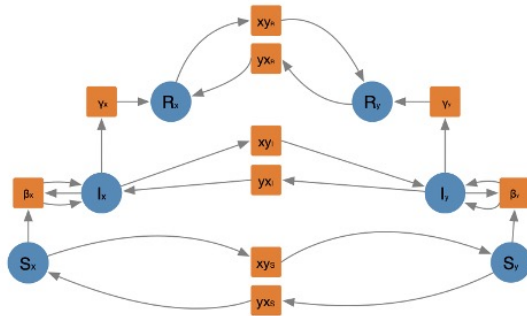
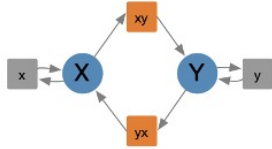
Sean Wu



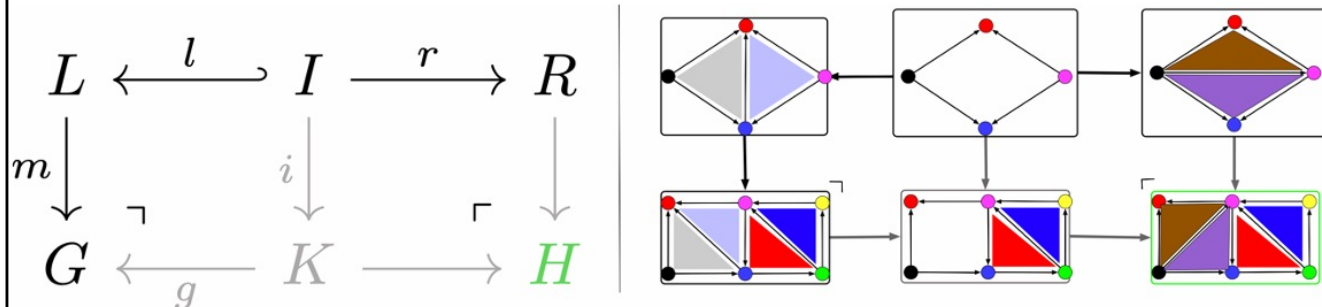
Andrew Baas



I'd like to thank all these people for welcoming me and being such fun collaborators!

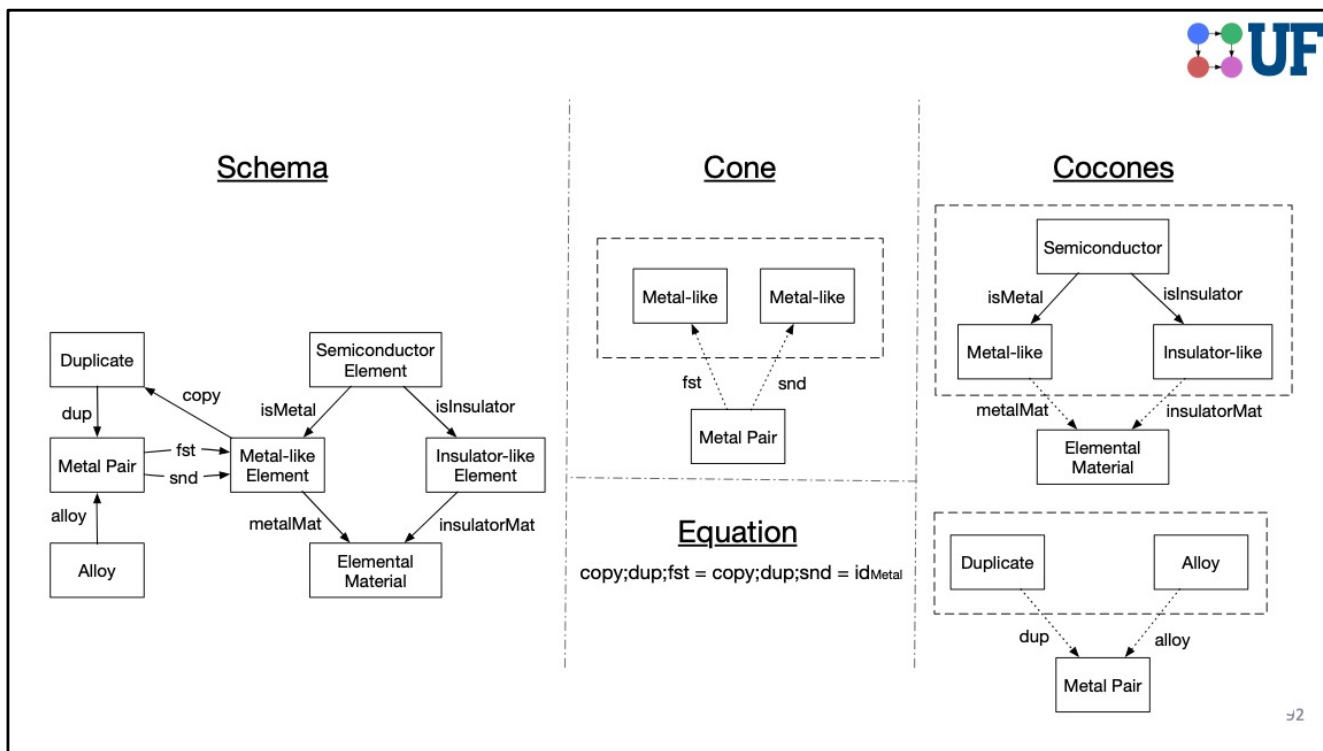


HIDDEN
 The REAL pullback (slightly more annoying to specify)



HIDDEN

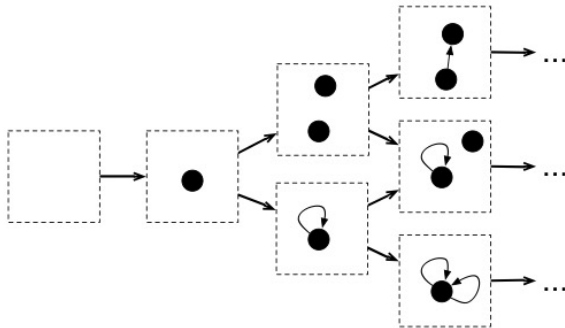
Here's another DPO example using the C-set of triangles I showed near the beginning of the talk, where we flip a quadrilateral.



HIDDEN

Example of scientific knowledge encoded by finite (co)limit sketches

II. Defining model space: enumeration



Semigroup order	1	2	3	4
# of C-sets	1	16	1968	429496729
# of semigroups	1	5	24	188

- If constraints were arbitrary code, we'd have no choice but to generate C-sets and filter.
- With a restricted language of constraints called *sketches*, we can greatly prune the search space.

- Pro: finds “unknown unknowns”
- Con: massive search space, results have no semantic meaning

93

HIDDEN

In some sense it's the most general but least interpretable way to explore models. I think sometimes we do want this - it can find structures we weren't consciously expecting.

And you might learn there's something wrong with your framework if this procedure produces models that seem fundamentally wrong. I'll show an example of that soon.

But the search space is huge, and unlike the other base-level ways of introducing generators, when your model has contents that get ultimately traced back to originating from enumeration, there's not much you can say about them, semantically.

Computationally, things were actually pretty easy when enumerating graphs because every valid C-set was a valid graph.

What about something with just a bit more structure, like a semigroup (which is an associative binary function).

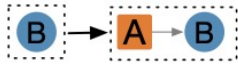
It turns out that a naïve enumeration algorithm will be useless for finding semigroups of a reasonable size.

So we want to constrain our search by something that has some structure that we can use to efficiently prune the space. It turns out there's a category-theoretic notion of a “sketch” which is a very expressive constraint language while still being computationally tractable.

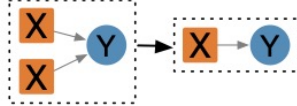
II. Defining model space: constraints

- Subset of first-order logic captured by C-Set homomorphisms

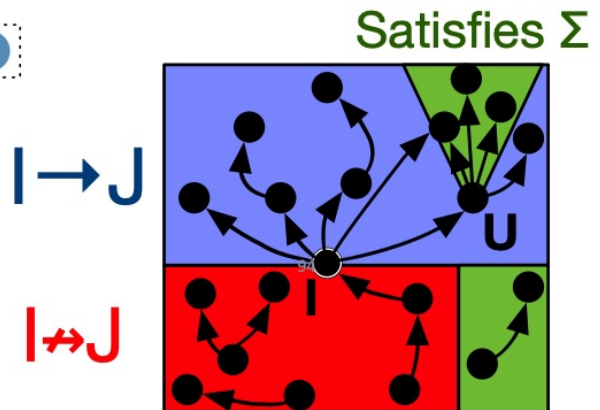
Outputs are surjective



Outputs are injective



- If I does not satisfy constraints Σ , we can algorithmically compute, U, the *best possible* relative of I, where:
 - I has been modified as little as possible
 - U still preserves the structure of I
 - U satisfies Σ



HIDDEN

Example of some constraints encoded by C-set homomorphisms

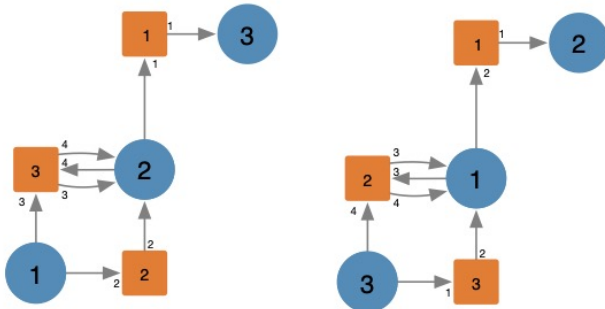
For example, if I want to say the output relation of a Petri net is surjective, then I want to say, for every state "B", there exists a transition "A" that has "B" as its output. For injectivity, I want to say, for every two transitions which output to the same state, those two transitions are actually the same. Together these would allow us to enforce that each species has a unique transition that produces it, if we wanted to do that for some reason.

You're specifying this complicated constraint specification and repair process, not with 1000 lines of opaque code, but just elegantly with these few structures.

III. Exploration of model space: note on symmetries



- We can arrive at the ‘same’ C-set via multiple paths
 - Rewriting rules, enumeration, products and sums, applying constraints
- # distinct C-sets on a given schema >>> # of isomorphism classes
- Distinctness due to implementation details: we want to ignore them



n	# C-Sets	# Iso classes	Time
1	1	1	Instant
2	256	17	Instant
3	531441	505	Minutes
4	4294967296	?	?

95

HIDDEN

it's possible to have models which we feel ought be considered to be equivalent actually seen as distinct. This makes our model space MUCH larger than it ought to be.

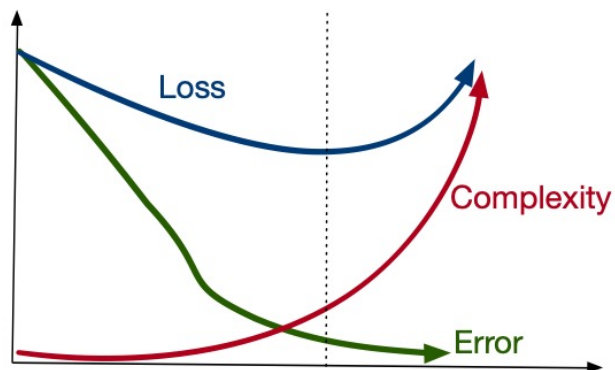
For example, to the computer, these two Petri Nets look completely different.

For a given pair, we can use a CSP algorithm to detect the existence of an isomorphism. But what if we have seen 1,000 models so far and want to check whether a freshly generated model is new? Rather than do 1000 pairwise tests, instead the solution is to just always work with a canonical permutation of the ACSet, which takes some effort to compute.

I have a blog post specifically about this problem and the generalization of a graph algorithm to ACSets, but the takeaway here is that when we explore models we are going to consider them only up to isomorphism.



Exploration: navigating a model space



- With monotonicity assumptions, we can find the optimal tradeoff of performance vs complexity without evaluating every model.

96

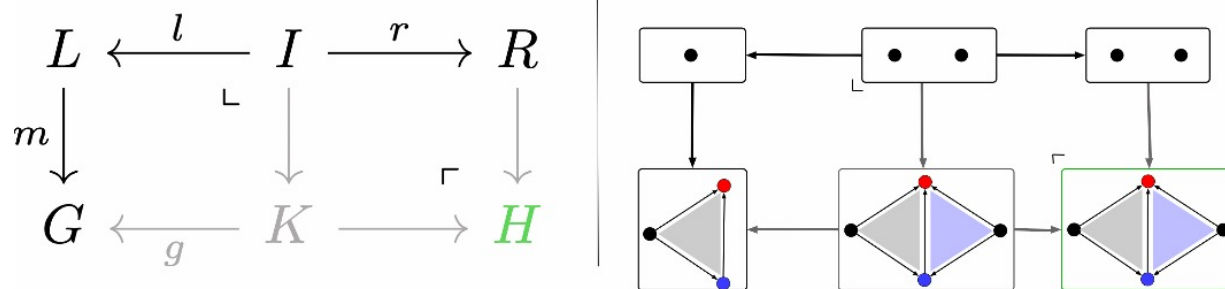
HIDDEN

If we pick our structure of our models and our loss function in a compatible way, then our loss profile will look smooth and predictable like this. That means we know when to stop exploring, once loss starts going up again.

This is a nice property to have, though not all loss functions will satisfy it.



Simulation: SqPO semantics



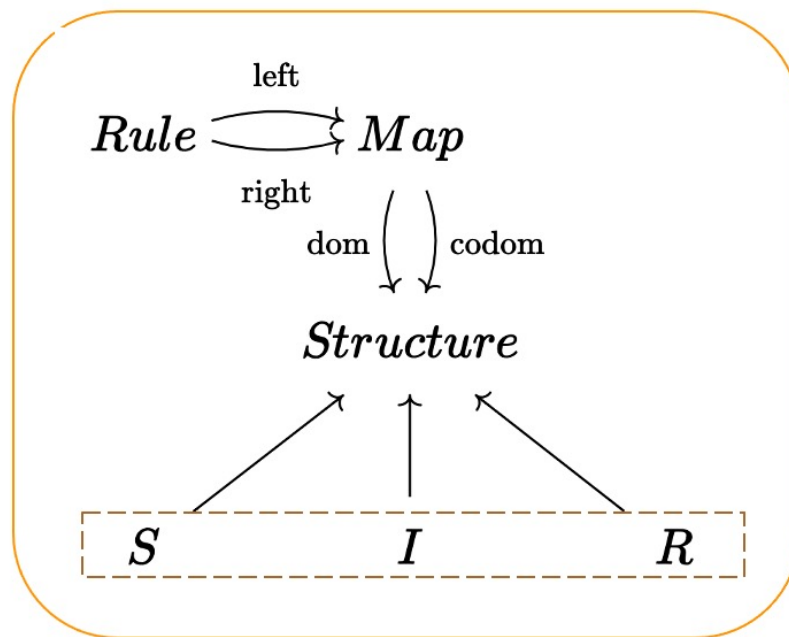
Brown et al. "Computational Category Theoretic Rewriting" (2022).

97

HIDDEN

Sesqui pushout rewriting adds a similar ability to SPO where you can add things which implicitly add other things.

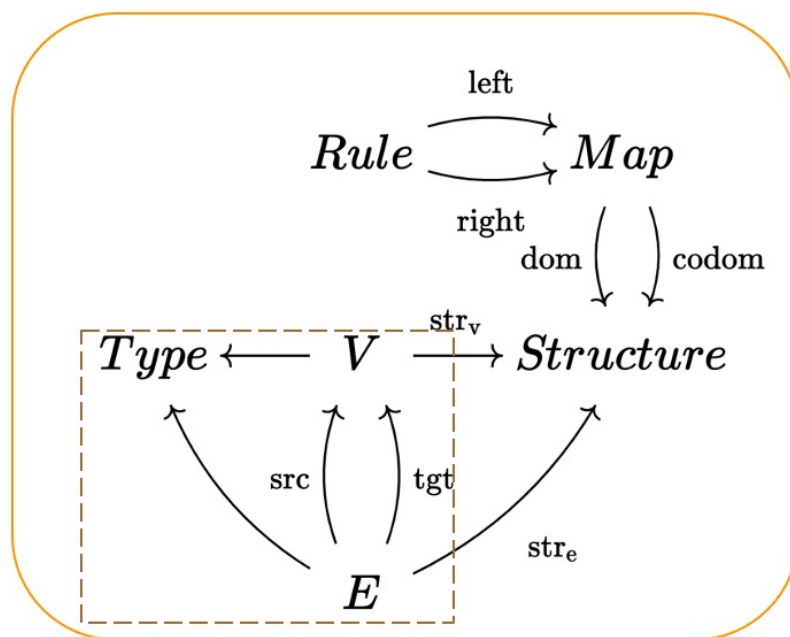
We can define one rewrite rule, which is a syntax, and freely give it any of these semantics.



HIDDEN

The data of a very basic agent-based model is captured in the following C-Set. It has multiple rewrite rules, each which has the data of a C-set homomorphism. Note we could have any type of C-set instead of this SIR C-Set by substituting a different C-set in the dotted box.

(We also want to add attributes to the rules to help us schedule them when running a simulation.)



HIDDEN

Here's an example of an agent based model with the state of world being represented with a typed graph.