

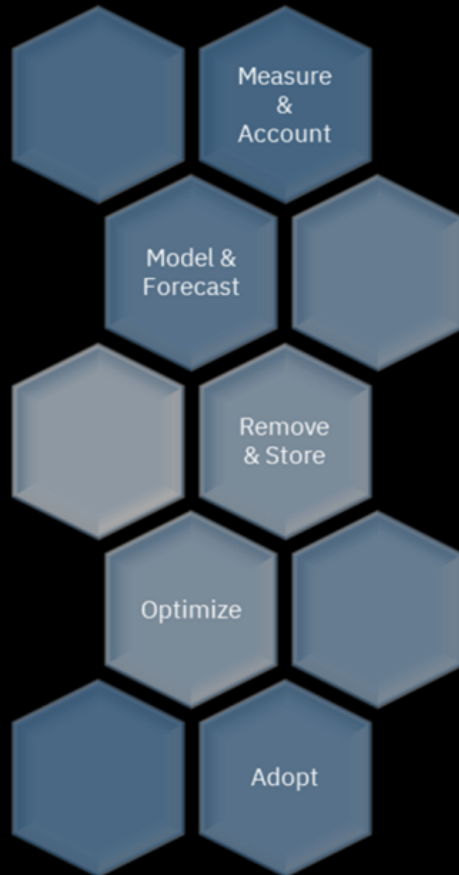
Harnessing the Geospatial Science Revolution  
to Advance Sustainability Research:

# A Geospatial Discovery Network

*C.D. Watson, H.F. Hamann, K. Weldemariam,  
B. Edwards, A. Jones, and J. Schmude*



# Accelerating computing for climate is critical for battling climate change



- *Greenhouse Gas Monitoring and Analysis*
- *Accelerating ESG reports*
- *Surrogate climate models*
- *AI for Materials Discovery*
  - *Carbon Capture*
  - *Membranes*
  - *Battery Materials*
- *Natural Carbon Sequestration*
- *Smarter transportation, buildings*
- *Next gen electric grid*
- *Climate impact modeling*
  - *Flood, wildfire, draughts, heatwaves etc.*
  - *Socio-economic*
- *Risk and impact discovery*
- ...

# Traditional systems are outpaced by explosion of data, complexity and urgency of the matter

## The power and promise of improved climate data infrastructure

Kevin Gurney<sup>a,1</sup> and Paul Shepson<sup>b</sup>

nation, and whether they're establishing the trust necessary to mobilize and sustain reduction investment.

### Problems with the Status Quo

Right now, however, US climate data collection and dissemination efforts are falling short. Measurement and tracking of GHG emissions reflect a collection of ad hoc mandates and voluntary interests. The measurement efforts are aimed at addressing a wide array of decision support needs with varying degrees of completeness and utility. For example, the US Environmental Protection Agency (EPA) produces the national GHG inventory as part of the ongoing obligation of the United States to the international treaty process (3). California has established a similar inventory of statewide emissions and is moving toward operational monitoring with surface

national. For example, knowing the total vehicular GHG emissions in a city or a state provides no direct insight into which roadways or which specific vehicle classes dominate the emissions or why. More specificity leads to greater efficiency and flexibility when, as is inevitable, we will have to choose which emissions to target first at lowest cost before tackling those that are smaller and more expensive to manage. Knowing who, where, and why emissions occur at local scales also assists in understanding emissions responsibility or ownership, which is key to any future policy that incorporates market mechanisms or trading of emission credits.

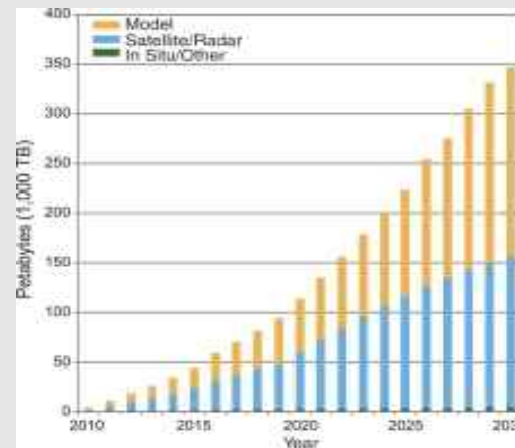
Furthermore, if emitters themselves are tasked with emissions measurement, the outcomes are open to internal bias, known to occur with "self-regulation" where emitters choose their emissions reduction target, decide which accounting methods to use, perform the

**Opinion: The power and promise of improved climate data infrastructure by Kevin Gurney, Paul Shepson**  
*Proceedings of the National Academy of Sciences* 118.35 (2021).

## Climate Data Challenges in the 21st Century

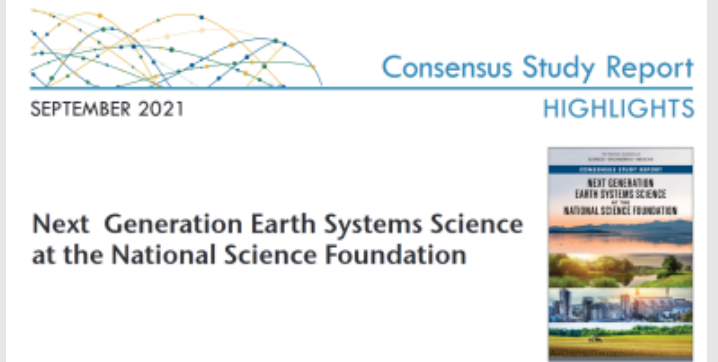
### Abstract

Climate data are dramatically increasing in volume and complexity, just as the users of these data in the scientific community and the public are rapidly increasing in number. **A new paradigm of more open, user-friendly data access is needed to ensure that society can reduce vulnerability to climate variability and change, while at the same time exploiting opportunities that will occur.**



### Climate Data Challenges in the 21st Century

By Jonathan T. Overpeck, Gerald A. Meehl, Sandrine Bony, David R. Easterling *Science* 11 700 (2001)



Use observational, computational, and modeling capabilities synergistically to accelerate discovery and convergence. The observational, computational, and modeling infrastructure must work collectively to support convergence in Earth Systems Science. Observations and monitoring reveal changes in the Earth's systems. Data from diverse sources are assimilated into models that represent natural- and social-system processes and their interactions across the Earth's systems. Computation provides the framework for putting together the complex pieces of Earth Systems Science, supporting data collection and analysis, generation of forecasts, and interpretation of model results.

# Key challenges are routed in the vast logical and physical distribution of climate-relevant information



The Amazon Sustainability Data Initiative

ASDI

Microsoft Planetary Computer



United States Census Bureau



Google Earth Engine



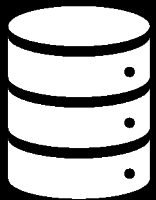
# Data gravity, lack of discoverability and modeling tools are amplifying the challenge



*Modeling frameworks*

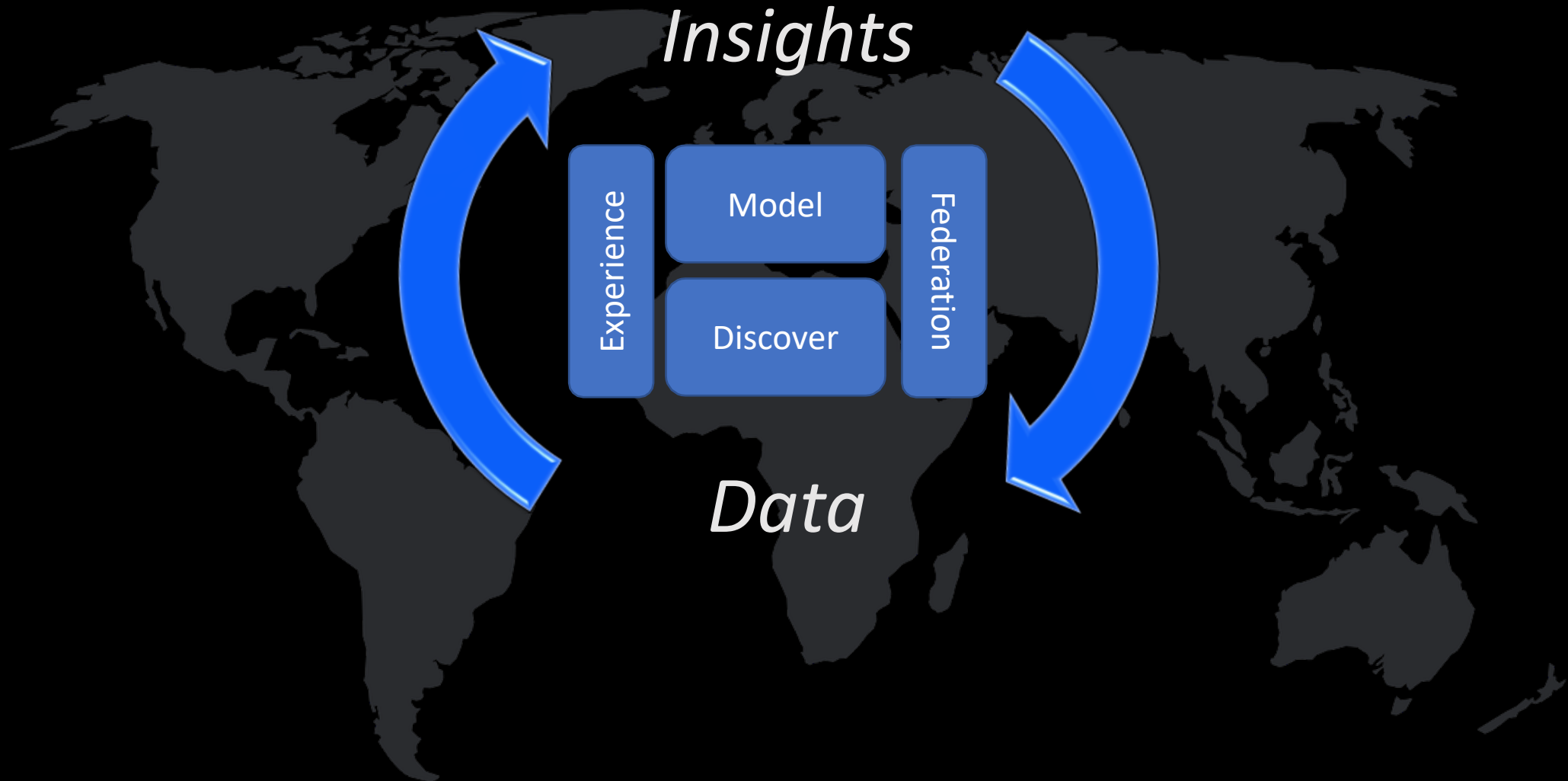


*Discoverability*

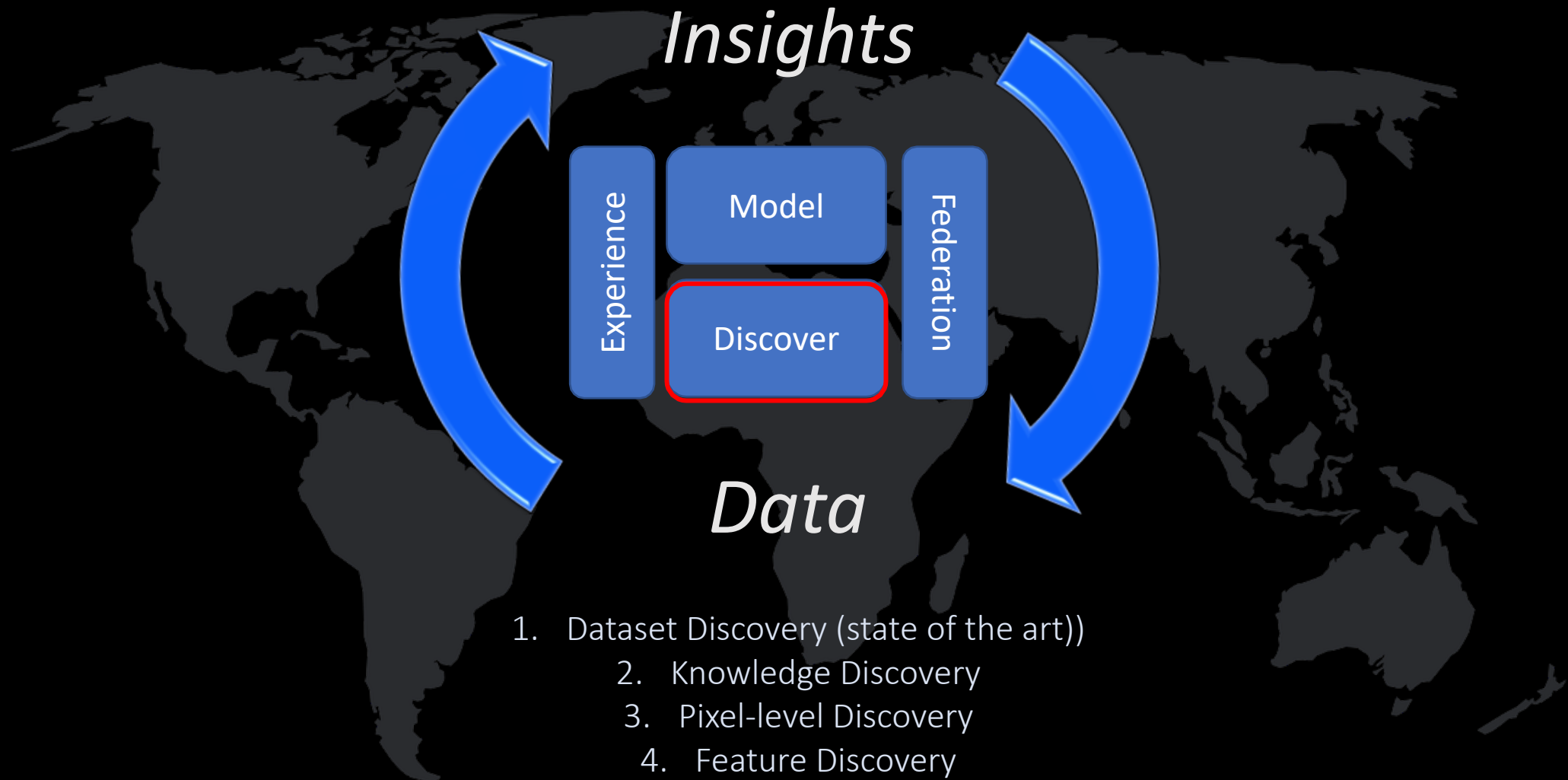


*Data Gravity*

# A Geospatial Discovery Network for distributed and federated discovery and modeling



# A Geospatial Discovery Network for distributed and federated discovery and modeling



# Accelerated Geospatial Knowledge Discovery with Deep Search

Document Conversion into Machine Readable Output

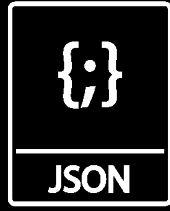
Climate Entity & Relationship Recognition

Climate Knowledge Representation

Geo-Spatial Document Visualization



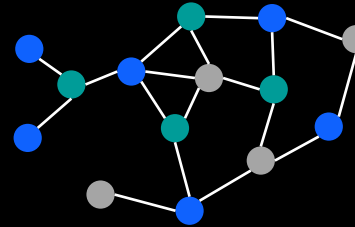
Convert



Extract



Load

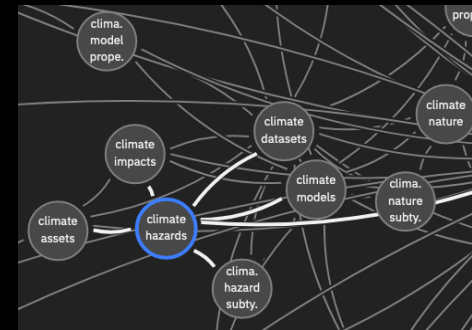
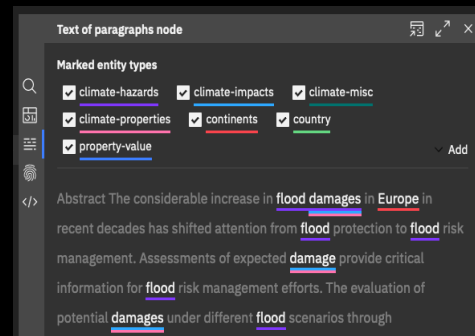


Query



## Knowledge Corpus

- 1 Mio. Climate related Abstracts
- IPCC Reports
- 8k ESG Reports
- 2 Mio. arXiv Publications



## Geospatial Knowledge Discovery

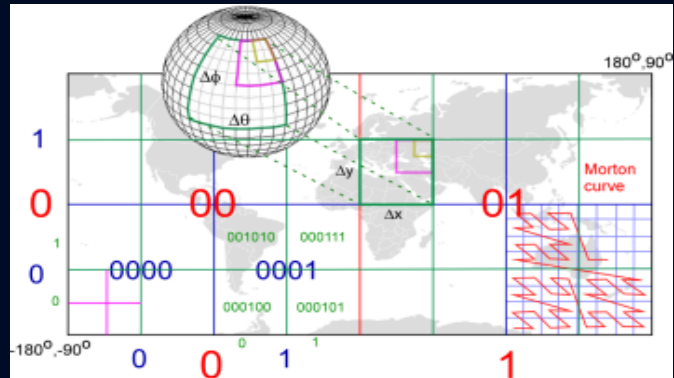
For a location and time:

- which surveys?
- what impact functions?
- what observations?
- what flood models?

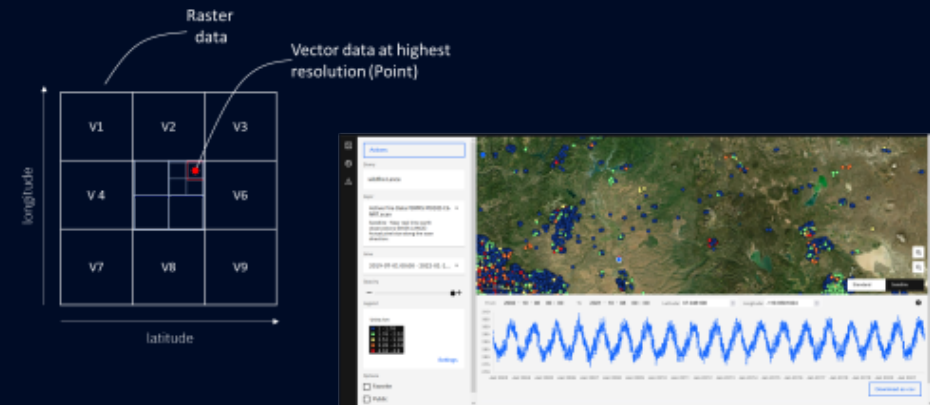


# Accelerated Pixel-level Discovery

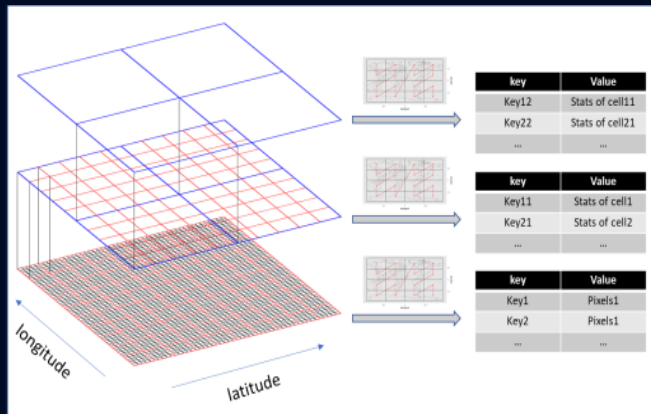
1. Nested resolution layers & common coordinate systems



3. Raster-Vector Data Cubes



2. Data/Pixel-level indexing with overview layers



⇒ Vector-Raster Data-level Discovery  
(e.g. find areas and timestamps with high population and precipitation rates 5x than climatology?)

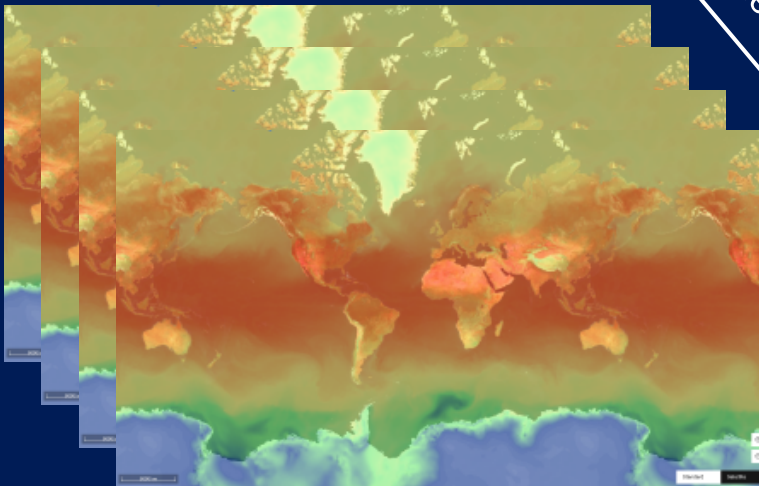


# Accelerated Feature Discovery

*Example: Spatially-resolved surface temperature gradients*

## 1. Available data

4.5M locations\* with 300,000 timestamps  
4.5M \* 4 Bytes \* 300,000 ~ 5.4 Terabytes per variable



## 2. User-defined functions

supporting arbitrary math to be submitted without downloads

```
def generate_query(layer_id,agg,startdate,duration,N,sql_name,poly):
    """Calculate the average annual value"""
    query_json = {
        'layers': [
            {
                'alias': '%(N)s'.format(str(years-startdate.year).zfill(2)),
                'aggregation': agg,
                'type': ' raster',
                'id': layer_id,
                'temporal': {'intervals': [
                    {'start': (startdate.replace(startdate.year + years)).strftime('%Y-%m-%d'),
                     'end': (startdate.replace(startdate.year + years) + timedelta(days = duration)).strftime('%Y-%m-%d')}
                ]},
                'output': False
            }
            for years in range(N,N)
        ],
        'spatial': {'type': 'poly', 'sql': sql, 'aggregation': ['sql'].poly},
        'temporal': {'intervals': [{}]},
        'name': name
    }

    """Calculate the mean of x"""
    query_json['layers'].extend([
        {
            'alias': 'm',
            'expression': '%(x)s'.join([
                '%(N)s'.format(str(years-startdate.year).zfill(2))
                for years in range(0,N)
            ])
        }
    ])

    query_json['layers'][N]['expression'] = query_json['layers'][N]['expression'] + '%(x)s' % str(x)
    query_json['layers'][N].update({'output': False})

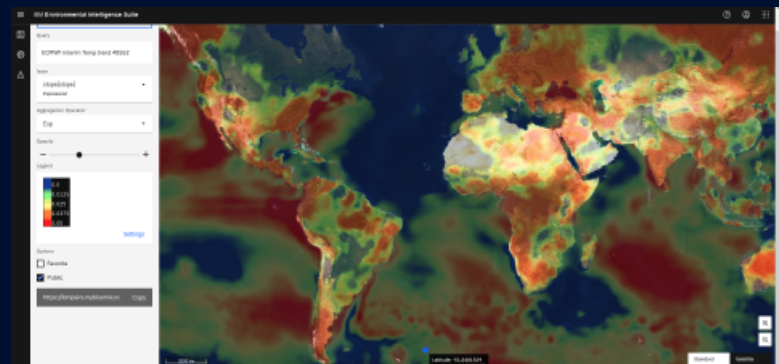
    """Calculate the mean of y"""
    query_json['layers'].extend([
        {
            'alias': 'y',
            'expression': '%(y)s'.join([
                '%(N)s'.format(str(years-startdate.year).zfill(2))
                for years in range(0,N)
            ])
        }
    ])

    query_json['layers'][N]['expression'] = query_json['layers'][N]['expression'] + '%(y)s' % str(y)
    query_json['layers'][N].update({'output': False})
```

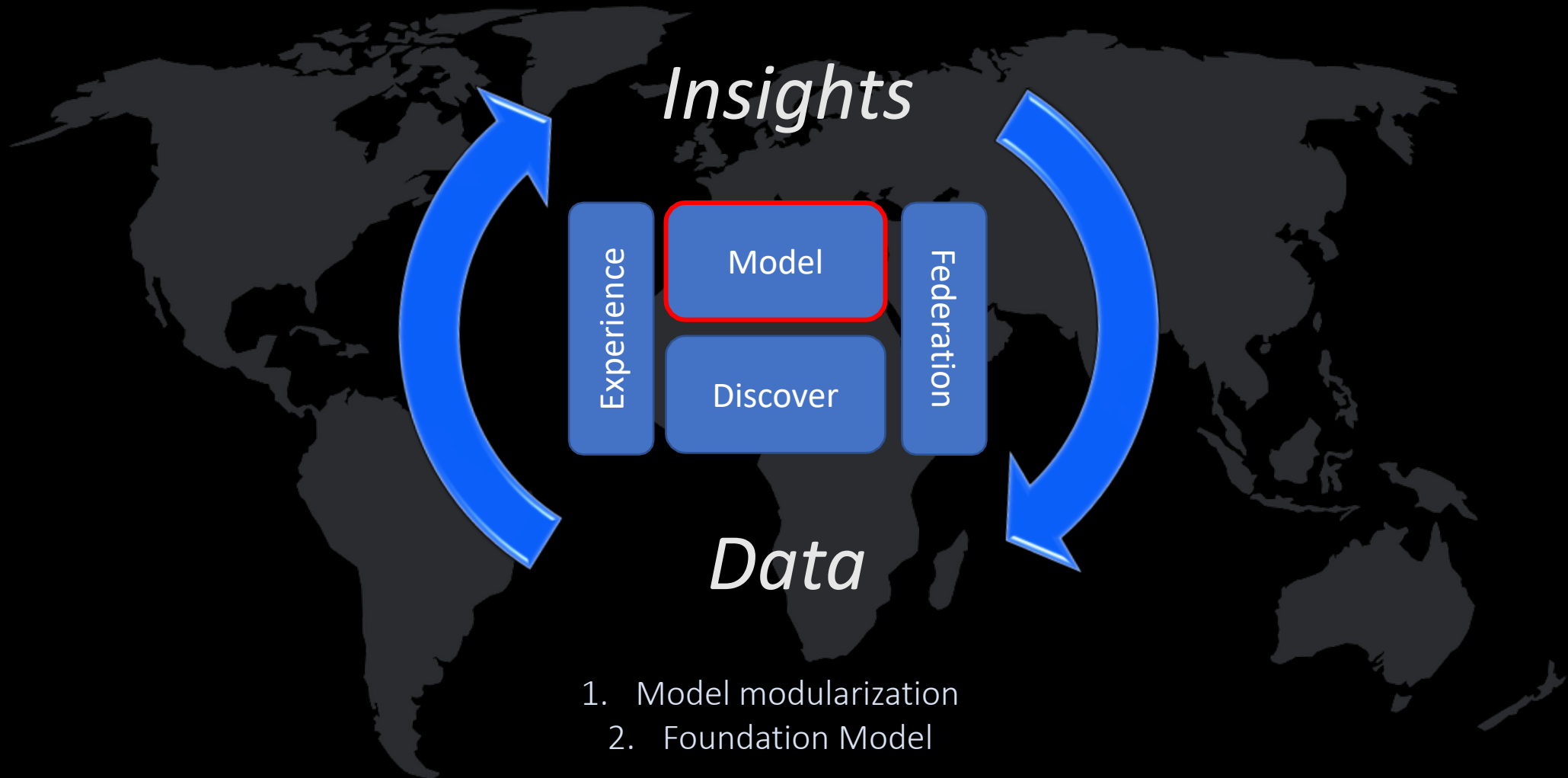
Time-series analysis (regression) over 300k timestamps for 4.5M locations

## 3. Features created

in 20 mins on a single thread

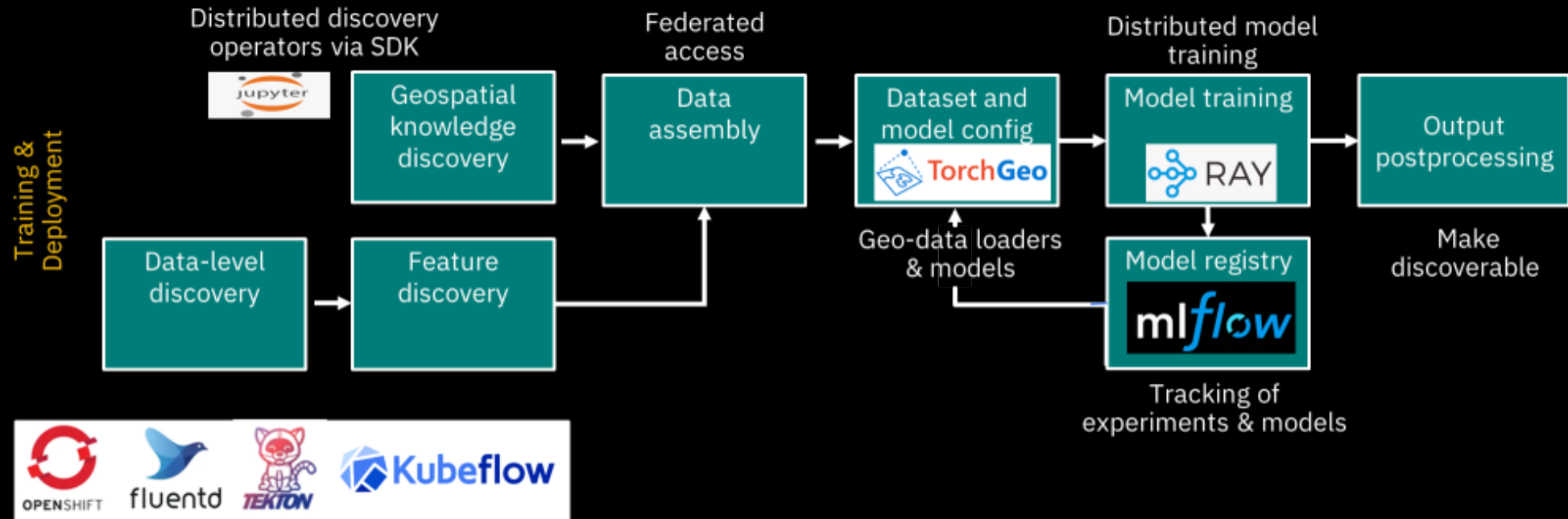


# A Geospatial Discovery Network for distributed and federated discovery and modeling



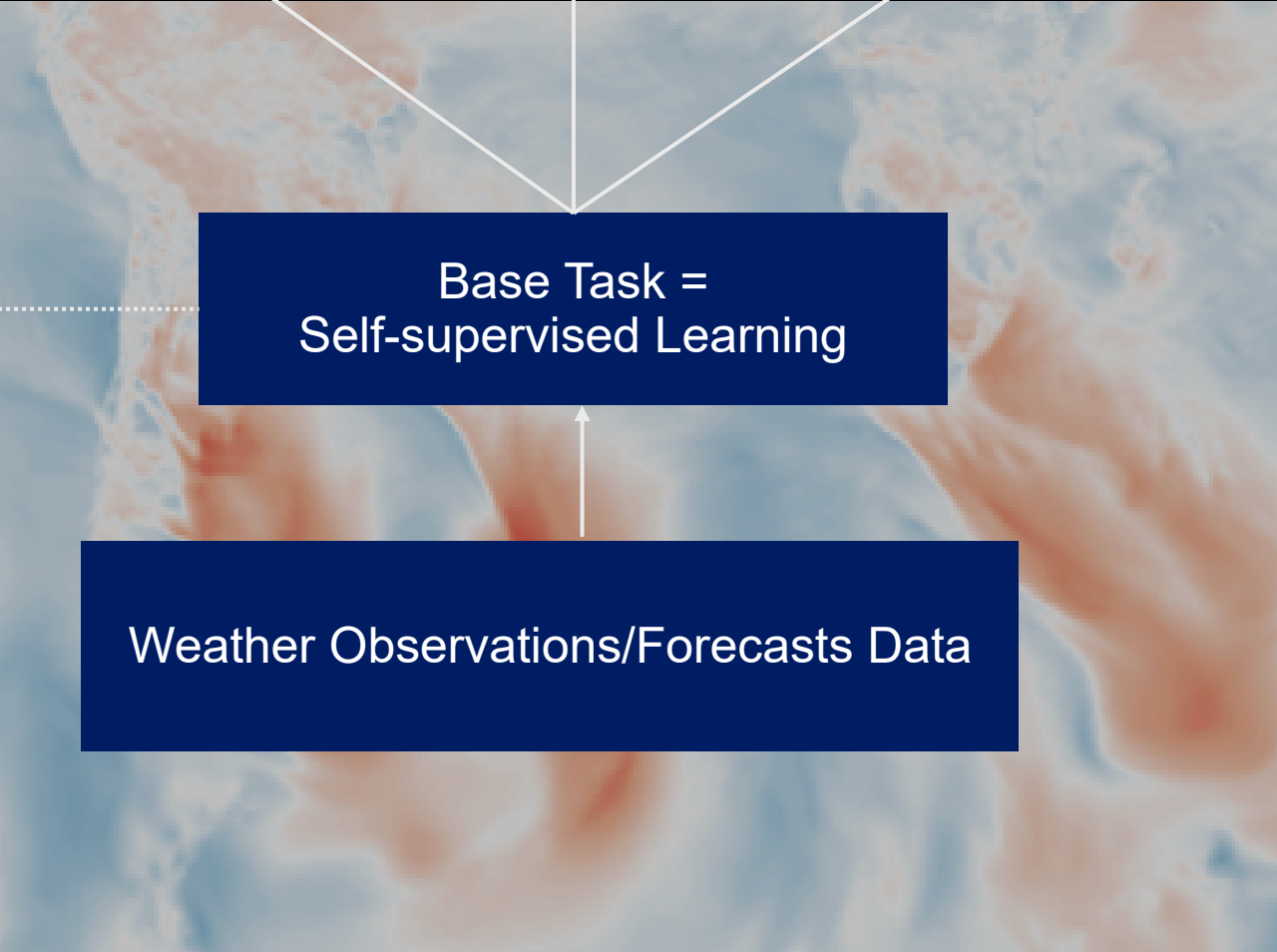
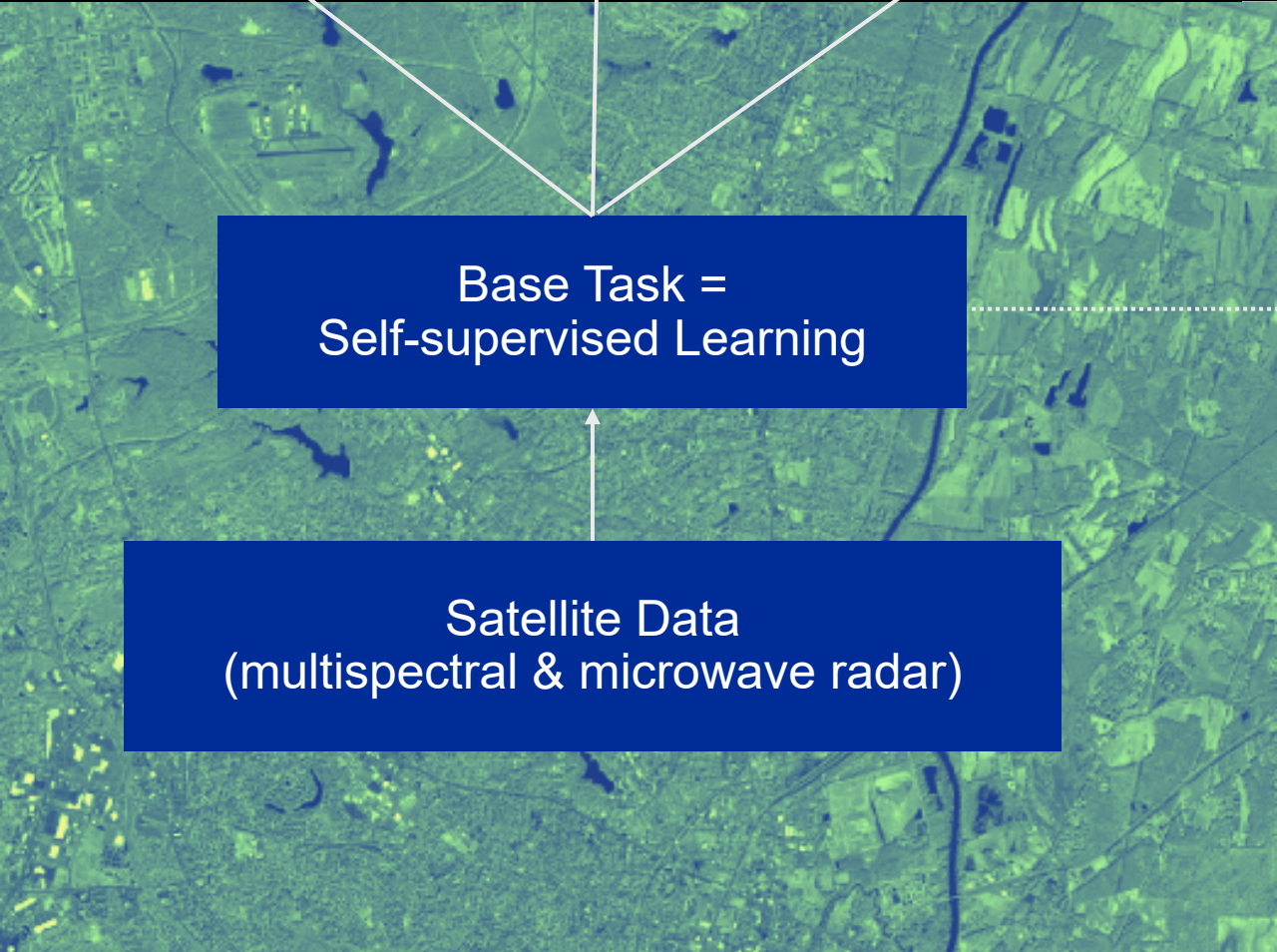
# Modular Modeling Framework

The screenshot displays a workflow management interface for an experiment named 'ifm-standard-ivncs'. On the left, a sidebar lists navigation options: Pipelines, Experiments, Runs, Recurring Runs, Artifacts, Executions, Documentation, and Github Repo. The main area shows a 'Graph' view of the pipeline with steps: 'static-pairs-query', 'precipitation-pairs-query', 'ifm-data-prep', 'ifm-precip-prep', and 'ifm'. A 'Simplify Graph' toggle is visible. On the right, a log window for a pod shows a series of messages including PAIRS base URL reads, JSON initialization, directory creation, data source opening, and successful query submissions.

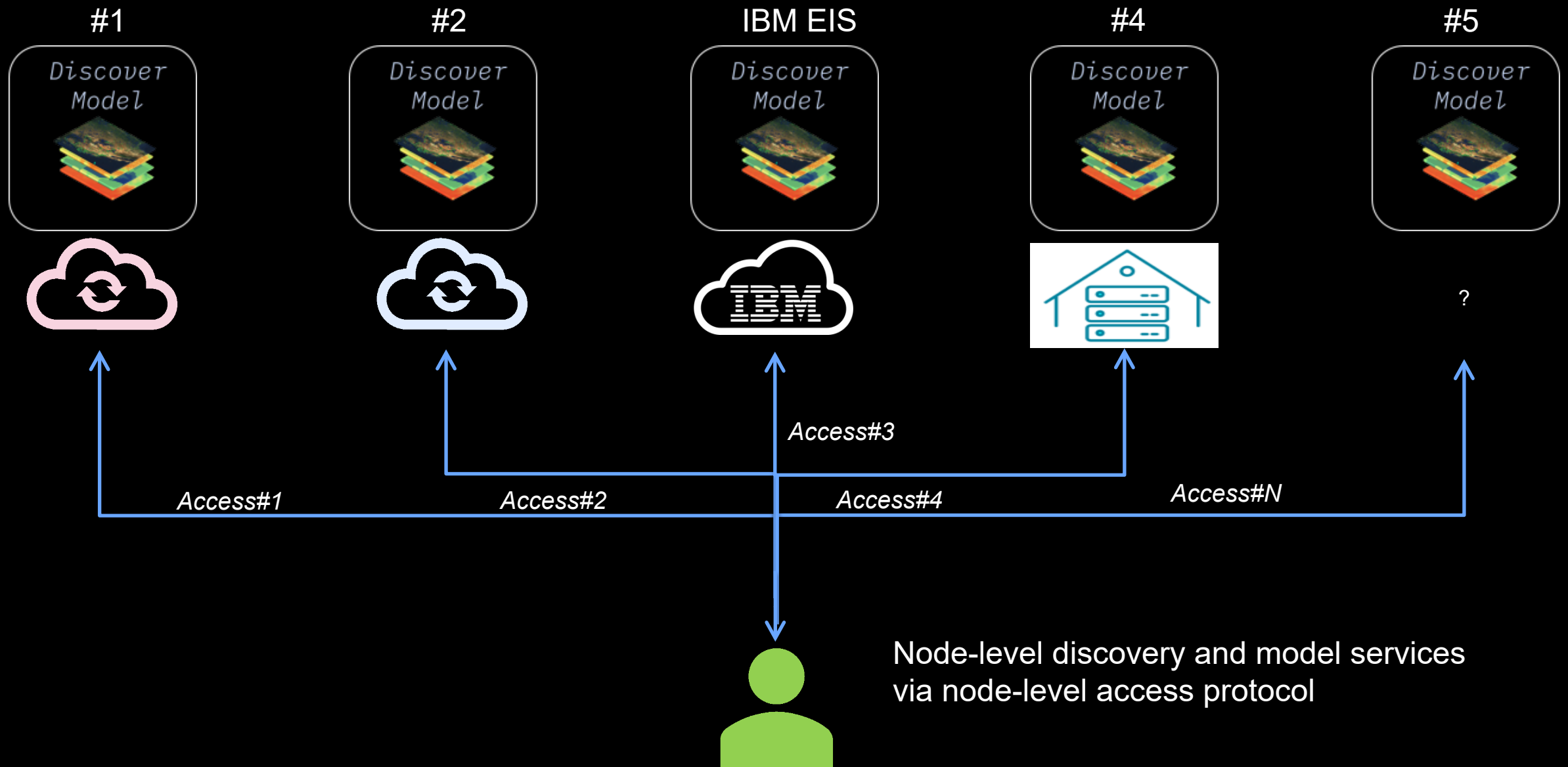


# Geospatial Foundation Models

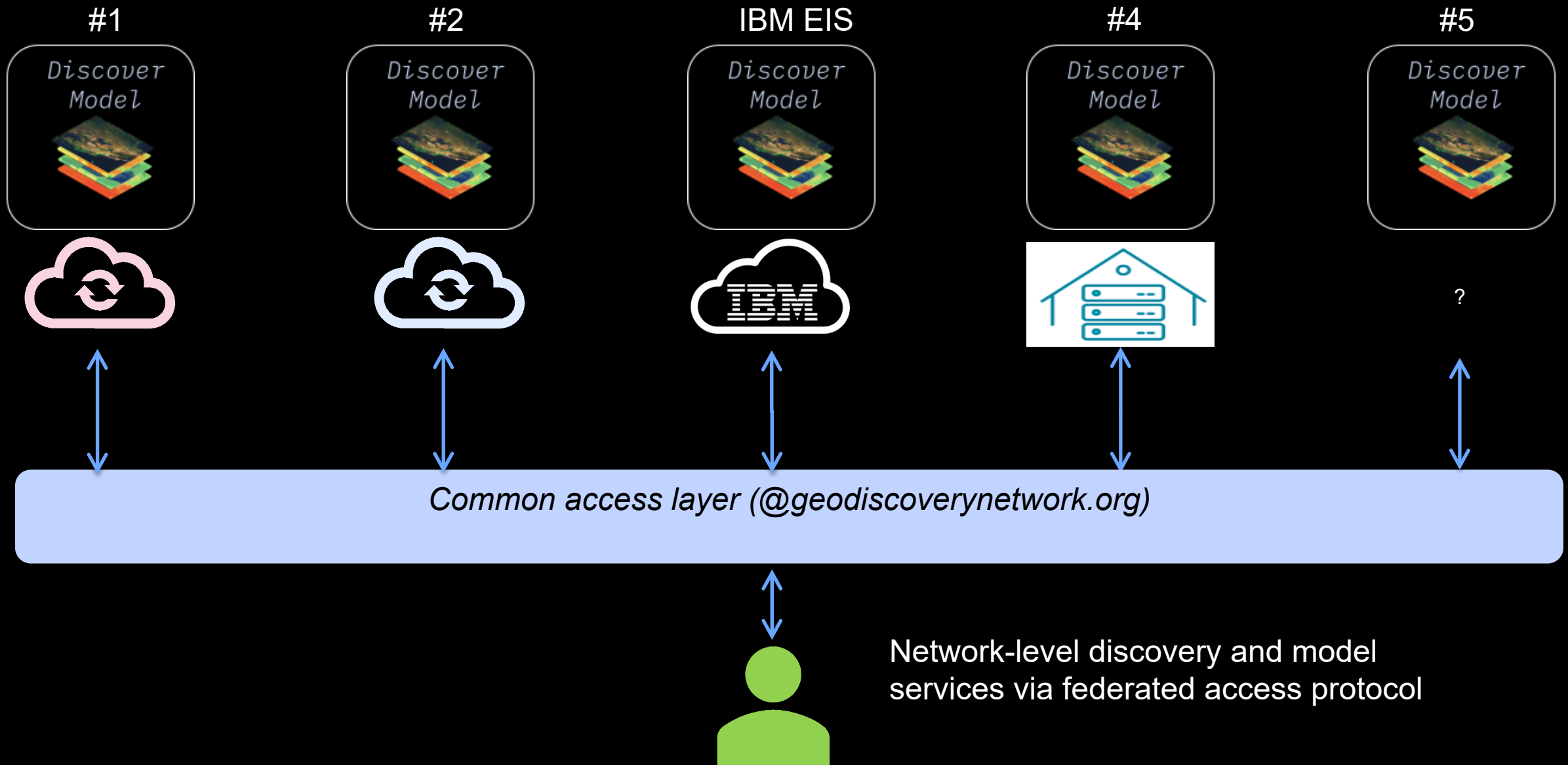
## Downstream Tasks



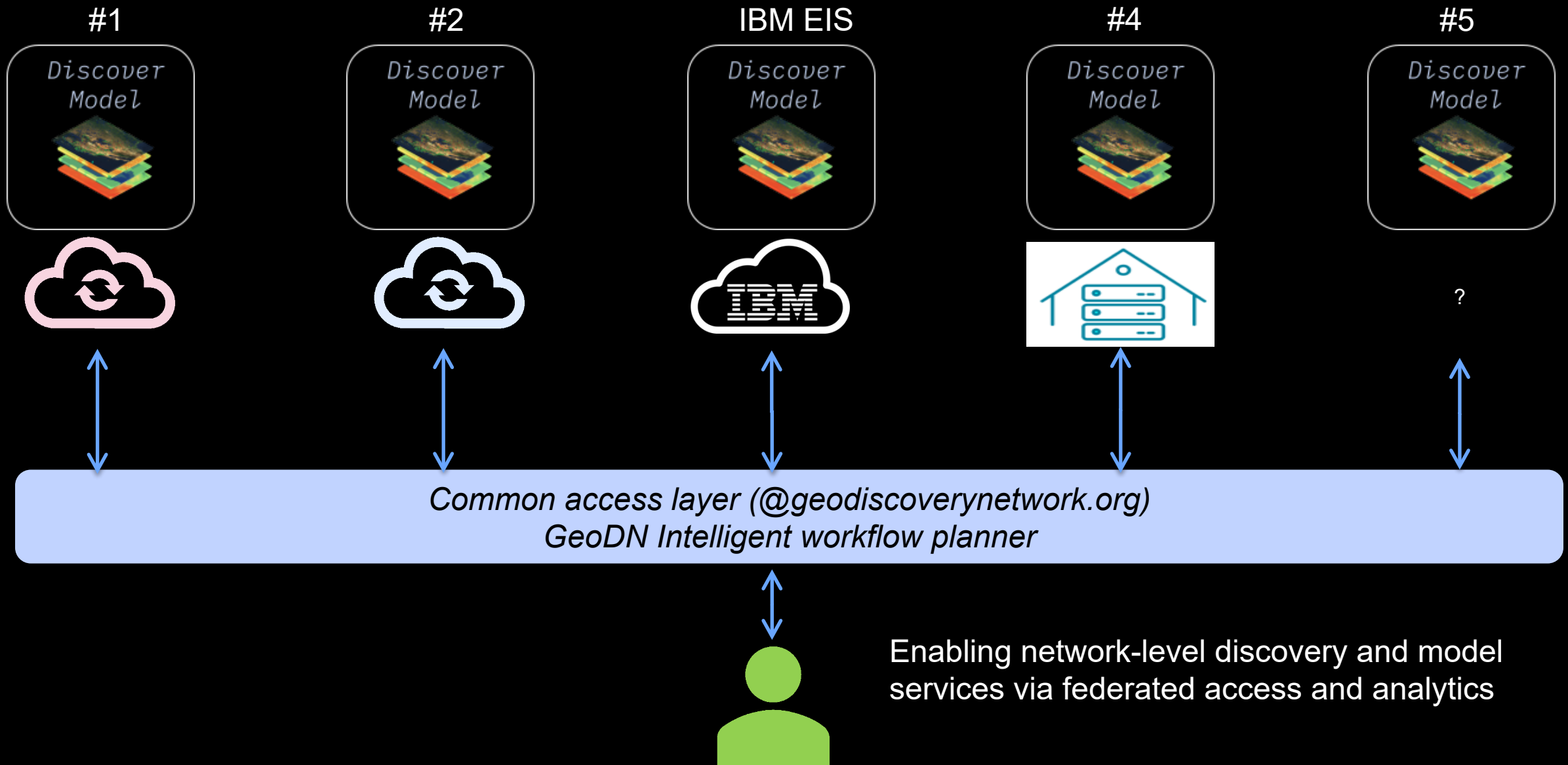
# Geospatial Discovery Network –Phase 1



# Geospatial Discovery Network –Phase 2



# Geospatial Discovery Network –Phase 3





Thank you

