



AutoMATES

Automated Model Assembly from Text, Equations and Software

Clayton Morrison

ml4ai.github.io/automates



ASKE Kickoff PI Meeting
5-6 December 2018

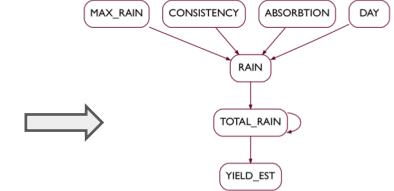
Motivation

Many detailed expert models encoded in software and described in text
... however...

- Software requires manual curation to integrate
- Textual descriptions and software are not integrated

Goal: Construct and curate semantically-rich representations of scientific models by integrating:
natural language descriptions and *equations* from publications and documentation, with the *software* that implements those models.

```
1 ====== UPDATE_EST - Updates the estimated yield of magic beans given
2 + rainfall and crop yield
3 + ======
4 +
5 + VARIABLES
6 +
7 + INPUT RAIN    = Additional rainfall
8 +
9 + ENDOUT YIELD_EST = Crop yield to update
10 +
11 +
12 ======
13 SUBROUTINE UPDATE_EST(MAIN, TOTAL_RAIN, YIELD_EST)
14   !-----+
15   !-----+ MAIN = Yield estimate, YIELD_EST = Total_Rain
16   !-----+ TOTAL_RAIN = Total Rain
17   !
18   !-----+
19   !-----+ Yield increases up to a point
20   !-----+
21   !-----+ IF(TOTAL_RAIN <= 40) THEN
22   !-----+
23   !-----+   YIELD_EST = -(TOTAL_RAIN - 40) ** 2 / 16 + 100
24   !-----+
25   !-----+ ELSEIF(TOTAL_RAIN > 100) THEN
26   !-----+
27   !-----+   YIELD_EST = -(TOTAL_RAIN - 100) ** 2 / 16 + 100
28   !-----+
29   !-----+ ELSE
30   !-----+
31   !-----+   YIELD_EST = -(TOTAL_RAIN - 40) ** 2 / 16 + 100
32   !-----+
33   !-----+ ENDIF
34   !-----+
35 END SUBROUTINE UPDATE_EST
36 .
37 .
```



Text

Subroutine LAIS is called for both phases to compute the change in leaf area index (dLAI). During vegetative period, LAI increases as a function of the rate of leaf number increase. The potential rate is limited by soil water stress (both deficit and saturation), through SWFAC, and temperature, through PT. Its value is given by:

<eqn 1>

Where PD is the plant density (plants/m²), EMP1 is the maximum leaf area expansion per leaf, (0.104 m²/leaf) and a is given by:

<eqn 2>

Where EMP2 and nb are coefficients in the expolinear equation and N is the development age of the plant (leaf number).

Equations

$$dLAI = SWFAC \cdot PT \cdot PD \cdot EMP1 \cdot \frac{a}{1+a} \quad <\text{eqn 1}>$$

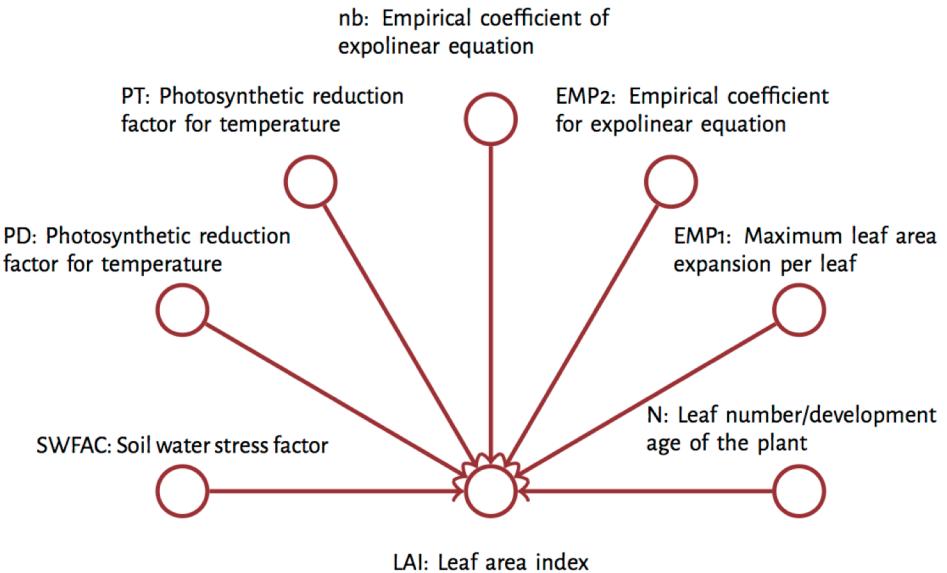
$$a = e^{EMP2 \cdot (N-nb)} \quad <\text{eqn 2}>$$

Software

- * dLAI = daily increase in leaf area index (m²/m²/d)
- * PD = plant density m⁻²
- * EMP1 = empirical coef. for expolinear eq.
- * EMP2 = empirical coef. for expolinear eq.
- * nb = empirical coef. for expolinear eq.
- * dN = incremental leaf number
- * N = leaf number
- * PT = photosynthesis reduction factor for temp.

$$a = \exp(EMP2 * (N-nb))$$

$$dLAI = SWFAC * PD * EMP1 * PT * (a/(1+a)) * dN$$



Semantically-Enriched
Grounded Function Network (GrFN)

Text

Subroutine LAIS is called for both phases to compute the change in leaf area index (dLAI). During vegetative period, LAI increases as a function of the rate of leaf number increase. The potential rate is limited by soil water stress (both deficit and saturation), through SWFAC, and temperature, through PT. Its value is given by:

<eqn 1>

Where PD is the plant density (plants/m²), EMP1 is the maximum leaf area expansion per leaf, (0.104 m²/leaf) and a is given by:

<eqn 2>

Where EMP2 and nb are coefficients in the expolinear equation and N is the development age of the plant (leaf number).

Equations

$$dLAI = SWFAC \cdot PT \cdot PD \cdot EMP1 \cdot \frac{a}{1+a} \quad <\text{eqn 1}>$$

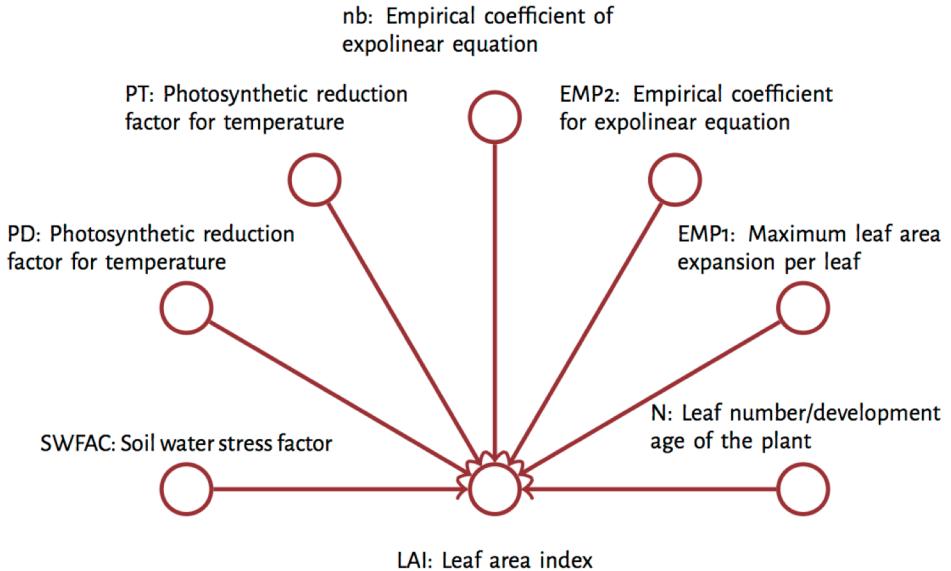
$$a = e^{EMP2 \cdot (N-nb)} \quad <\text{eqn 2}>$$

Software

- * dLAI = daily increase in leaf area index (m²/m²/d)
- * PD = plant density m⁻²
- * EMP1 = empirical coef. for expolinear eq.
- * EMP2 = empirical coef. for expolinear eq.
- * nb = empirical coef. for expolinear eq.
- * dN = incremental leaf number
- * N = leaf number
- * PT = photosynthesis reduction factor for temp.

$$a = \exp(EMP2 * (N-nb))$$

$$dLAI = SWFAC * PD * EMP1 * PT * (a/(1+a)) * dN$$



Semantically-Enriched
Grounded Function Network (GrFN)

Text

Subroutine LAIS is called for both phases to compute the change in leaf area index (dLAI). During vegetative period, LAI increases as a function of the rate of leaf number increase. The potential rate is limited by soil water stress (both deficit and saturation), through SWFAC, and temperature, through PT. Its value is given by:

<eqn 1>

Where PD is the plant density (plants/m²), EMP1 is the maximum leaf area expansion per leaf, (0.104 m²/leaf) and a is given by:

<eqn 2>

Where EMP2 and nb are coefficients in the expolinear equation and N is the development age of the plant (leaf number).

Equations

$$dLAI = SWFAC \cdot PT \cdot PD \cdot EMP1 \cdot \frac{a}{1+a} \quad <\text{eqn 1}>$$

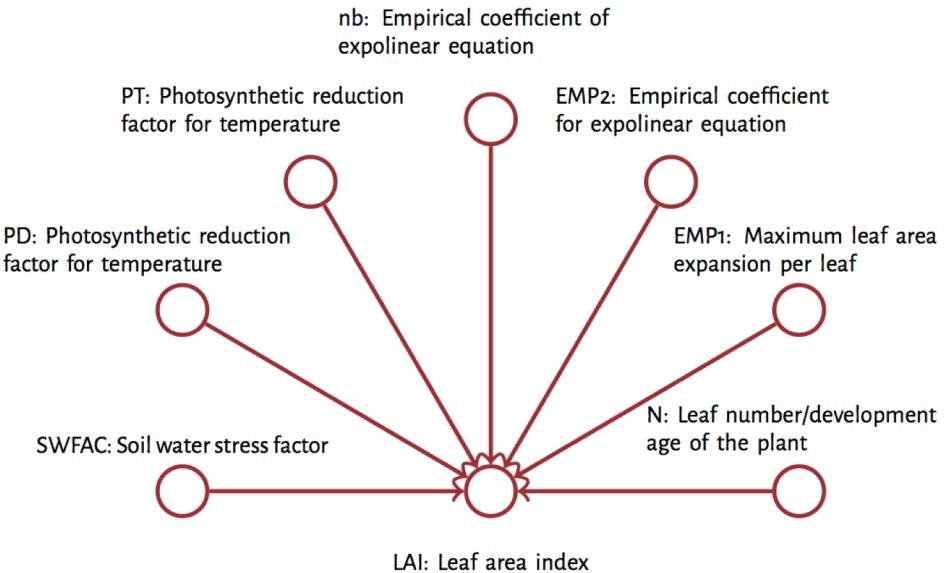
$$a = e^{EMP2 \cdot (N-nb)} \quad <\text{eqn 2}>$$

Software

- * dLAI = daily increase in leaf area index (m²/m²/d)
- * PD = plant density m⁻²
- * EMP1 = empirical coef. for expolinear eq.
- * EMP2 = empirical coef. for expolinear eq.
- * nb = empirical coef. for expolinear eq.
- * dN = incremental leaf number
- * N = leaf number
- * PT = photosynthesis reduction factor for temp.

$$a = \exp(EMP2 * (N-nb))$$

$$dLAI = SWFAC * PD * EMP1 * PT * (a/(1+a)) * dN$$



Semantically-Enriched
Grounded Function Network (GrFN)

Text

Subroutine LAIS is called for both phases to compute the change in leaf area index (dLAI). During vegetative period, LAI increases as a function of the rate of leaf number increase. The potential rate is limited by soil water stress (both deficit and saturation), through SWFAC, and temperature, through PT. Its value is given by:

<eqn 1>

Where PD is the plant density (plants/m²), EMP1 is the maximum leaf area expansion per leaf, (0.104 m²/leaf) and a is given by:

<eqn 2>

Where EMP2 and nb are coefficients in the expolinear equation and N is the development age of the plant (leaf number).

Equations

$$dLAI = SWFAC \cdot PT \cdot PD \cdot EMP1 \cdot \frac{a}{1+a} \quad <\text{eqn 1}>$$

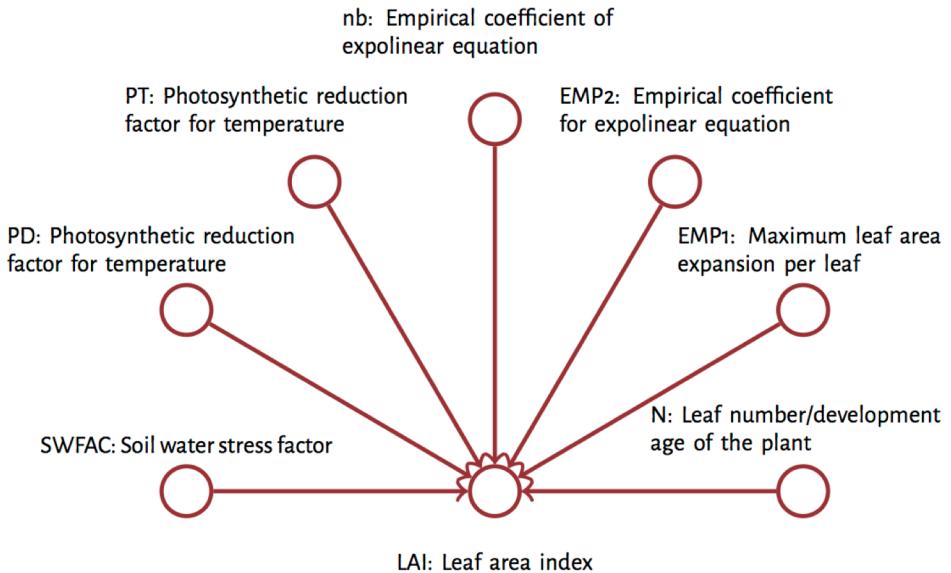
$$a = e^{EMP2 \cdot (N-nb)} \quad <\text{eqn 2}>$$

Software

- * dLAI = daily increase in leaf area index (m²/m²/d)
- * PD = plant density m⁻²
- * EMP1 = empirical coef. for expolinear eq.
- * EMP2 = empirical coef. for expolinear eq.
- * nb = empirical coef. for expolinear eq.
- * dN = incremental leaf number
- * N = leaf number
- * PT = photosynthesis reduction factor for temp.

$$a = \exp(EMP2 * (N-nb))$$

$$dLAI = SWFAC * PD * EMP1 * PT * (a/(1+a)) * dN$$



Semantically-Enriched
Grounded Function Network (GrFN)

Text

Subroutine LAIS is called for both phases to compute the change in leaf area index (dLAI). During vegetative period, LAI increases as a function of the rate of leaf number increase. The potential rate is limited by soil water stress (both deficit and saturation), through SWFAC, and temperature, through PT. Its value is given by: $\text{}$

Where PD is the plant density (plants/m²), EMP1 is the maximum leaf area expansion per leaf, (0.104 m²/leaf) and a is given by: $\text{}$

Where EMP2 and nb are coefficients in the expolinear equation and N is the development age of the plant (leaf number).

Equations

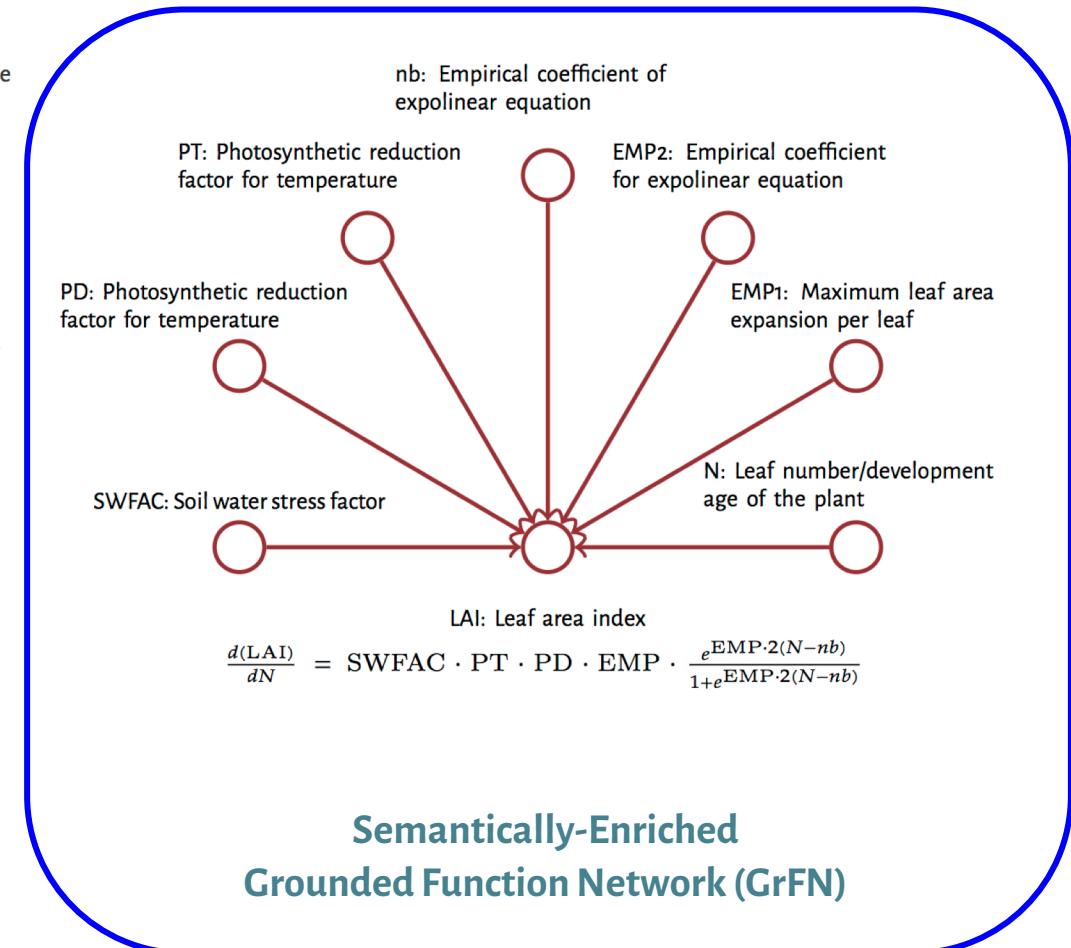
$$dLAI = SWFAC \cdot PT \cdot PD \cdot EMP1 \cdot \frac{a}{1+a} \quad \text{}$$
$$a = e^{EMP2 \cdot (N-nb)} \quad \text{}$$

Software

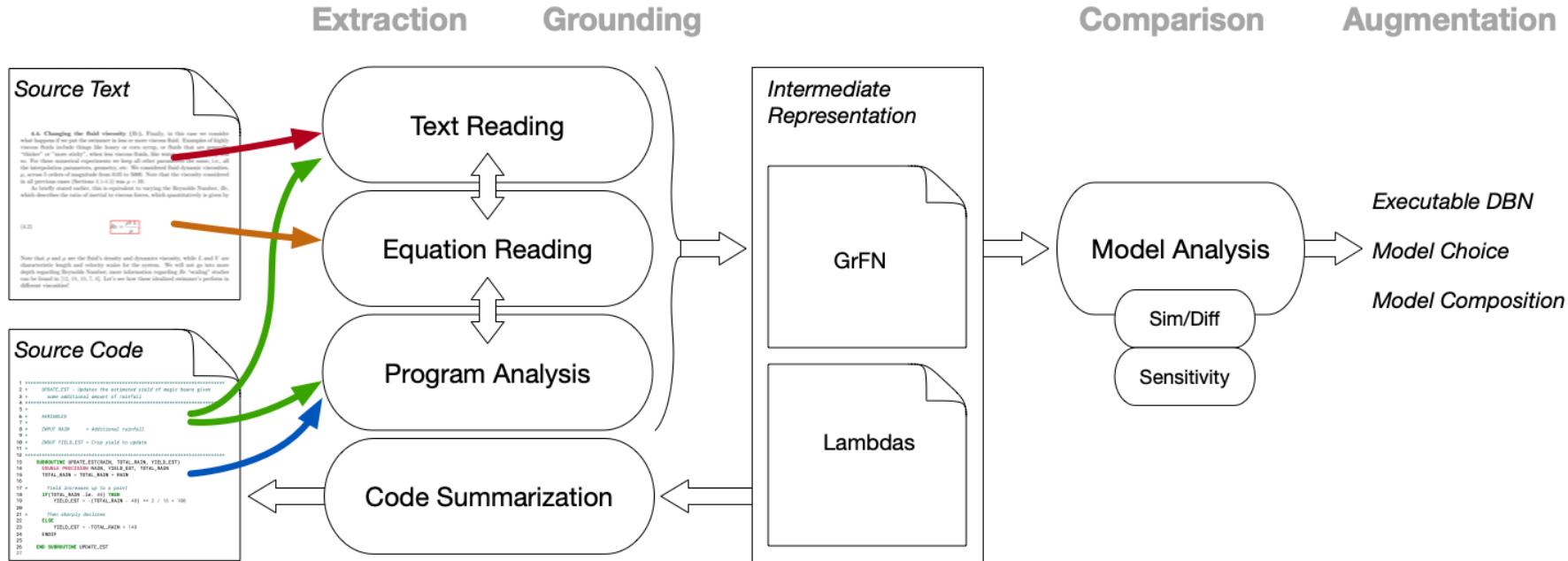
- * dLAI = daily increase in leaf area index (m²/m²/d)
- * PD = plant density m⁻²
- * EMP1 = empirical coef. for expolinear eq.
- * EMP2 = empirical coef. for expolinear eq.
- * nb = empirical coef. for expolinear eq.
- * dN = incremental leaf number
- * N = leaf number
- * PT = photosynthesis reduction factor for temp.

$$a = \exp(EMP2 * (N-nb))$$

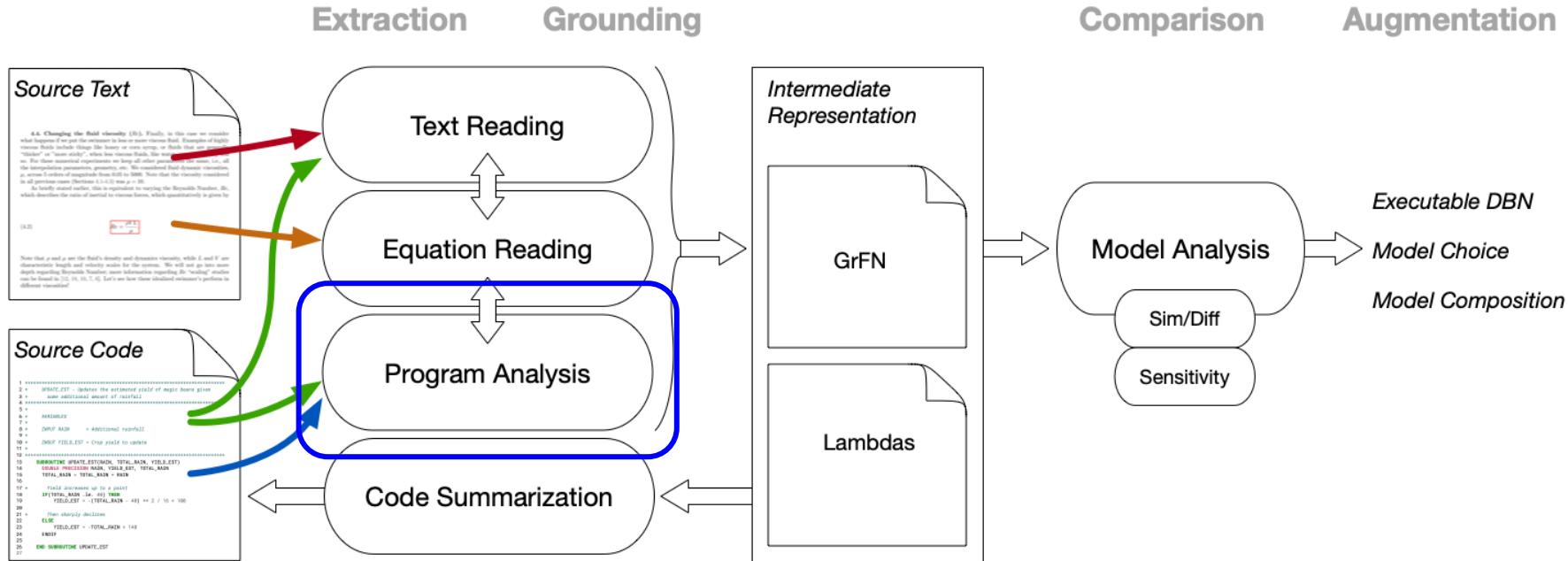
$$dLAI = SWFAC * PD * EMP1 * PT * (a/(1+a)) * dN$$



AutoMATES Architecture



AutoMATES Architecture



From source...

crop_yield.f

```
1 ****
2 *      UPDATE_EST - Updates the estimated yield of magic beans given
3 *      some additional amount of rainfall
4 ****
5 *
6 *      VARIABLES
7 *
8 *      INPUT RAIN      = Additional rainfall
9 *
10 *     INPUT YIELD_EST = Crop yield to update
11 *
12 ****
13 SUBROUTINE UPDATE_EST(RAIN, TOTAL_RAIN, YIELD_EST)
14
15     TOTAL_RAIN = TOTAL_RAIN + RAIN
16
17 *      Yield increases up to a point
18 IF(TOTAL_RAIN <= 40) THEN
19     YIELD_EST = -(TOTAL_RAIN - 40) ** 2 / 16 + 10
20
21 *      Then sharply declines
22 ELSE
23     YIELD_EST = -TOTAL_RAIN + 140
24 ENDIF
25
26 END SUBROUTINE UPDATE_EST
27
```

function
(program, subroutine)

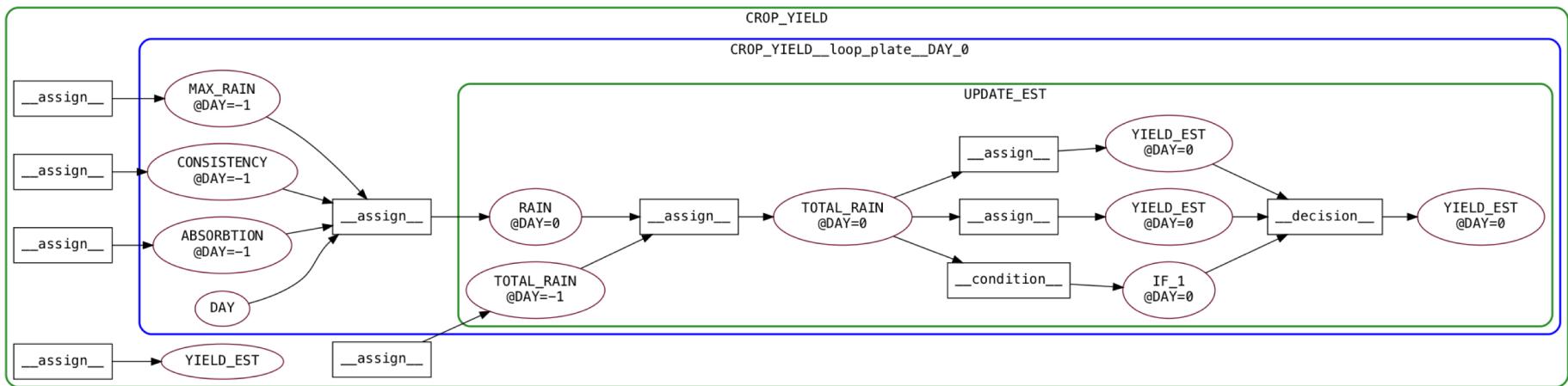
conditional

```
28 ****
29 *      CROP_YIELD - Estimate the yield of magic beans given a simple
30 *      model for rainfall
31 ****
32 *
33 *      VARIABLES
34 *
35 *      INPUT MAX_RAIN      = The maximum rain for the month
36 *      INPUT CONSISTENCY    = The consistency of the rainfall
37 *          (higher = more consistent)
38 *      INPUT ABSORBTION = Estimates the % of rainfall absorbed into the
39 *          soil (i.e. % lost due to evaporation, runoff)
40 *
41 *      OUTPUT YIELD_EST = The estimated yield of magic beans
42 *
43 *      DAY              = The current day of the month
44 *      RAIN             = The rainfall estimate for the current day
45 *
46 ****
47 PROGRAM CROP_YIELD
48 IMPLICIT NONE
49
50 INTEGER DAY
51 DOUBLE PRECISION RAIN, YIELD_EST, TOTAL_RAIN
52 DOUBLE PRECISION MAX_RAIN, CONSISTENCY, ABSORBTION
53
54 MAX_RAIN = 4.0
55 CONSISTENCY = 64.0
56 ABSORBTION = 0.6
57
58 YIELD_EST = 0
59 TOTAL_RAIN = 0
60
61 DO 20 DAY=1,31
62     Compute rainfall for the current day
63     RAIN = (-(DAY - 16) ** 2 / CONSISTENCY + MAX_RAIN) * ABSORBTION
64
65     Update rainfall estimate
66     CALL UPDATE_EST(RAIN, TOTAL_RAIN, YIELD_EST)
67     PRINT *, Day, Estimate:, YIELD_EST
68
69 20 ENDDO
70
71 PRINT *, "Crop Yield(%): ", YIELD_EST
72
73 END PROGRAM CROP_YIELD
```

var declarations

var assignments

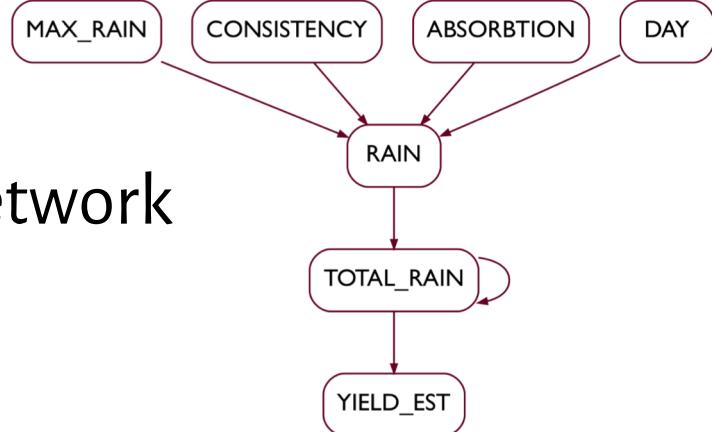
loop



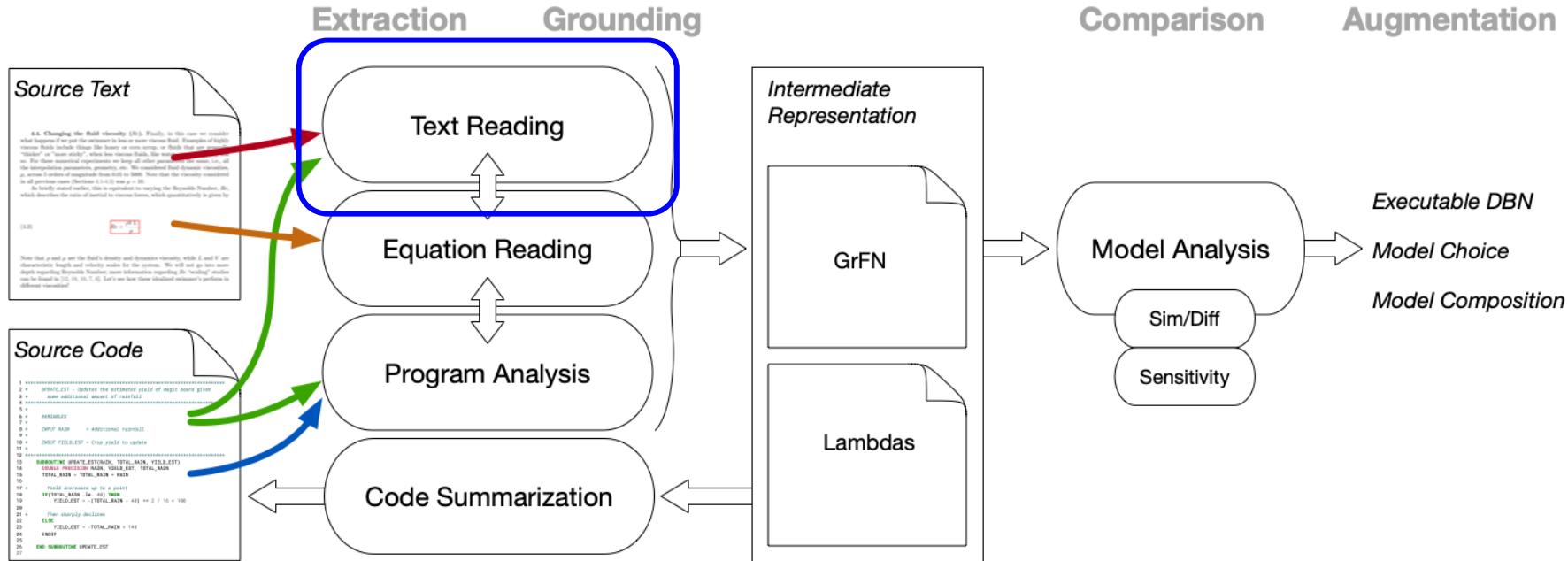
... derive:

GrFN: Grounded Function Network Executable DBN

... with associated comments

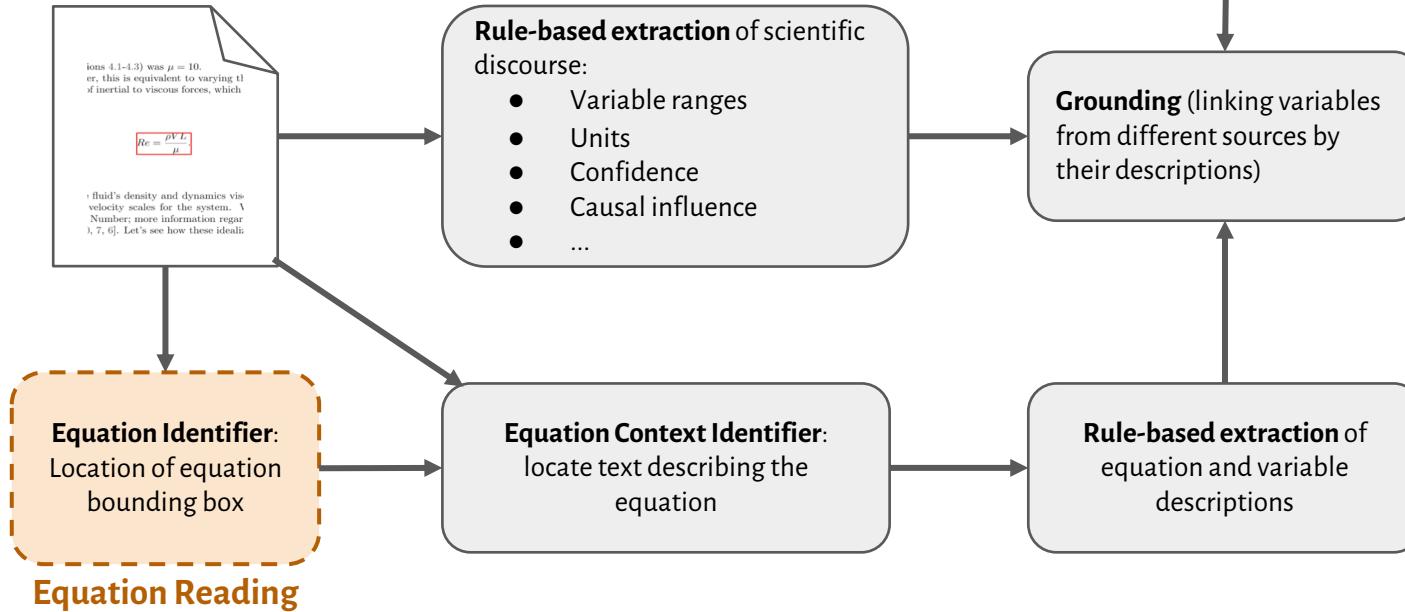


AutoMATES Architecture

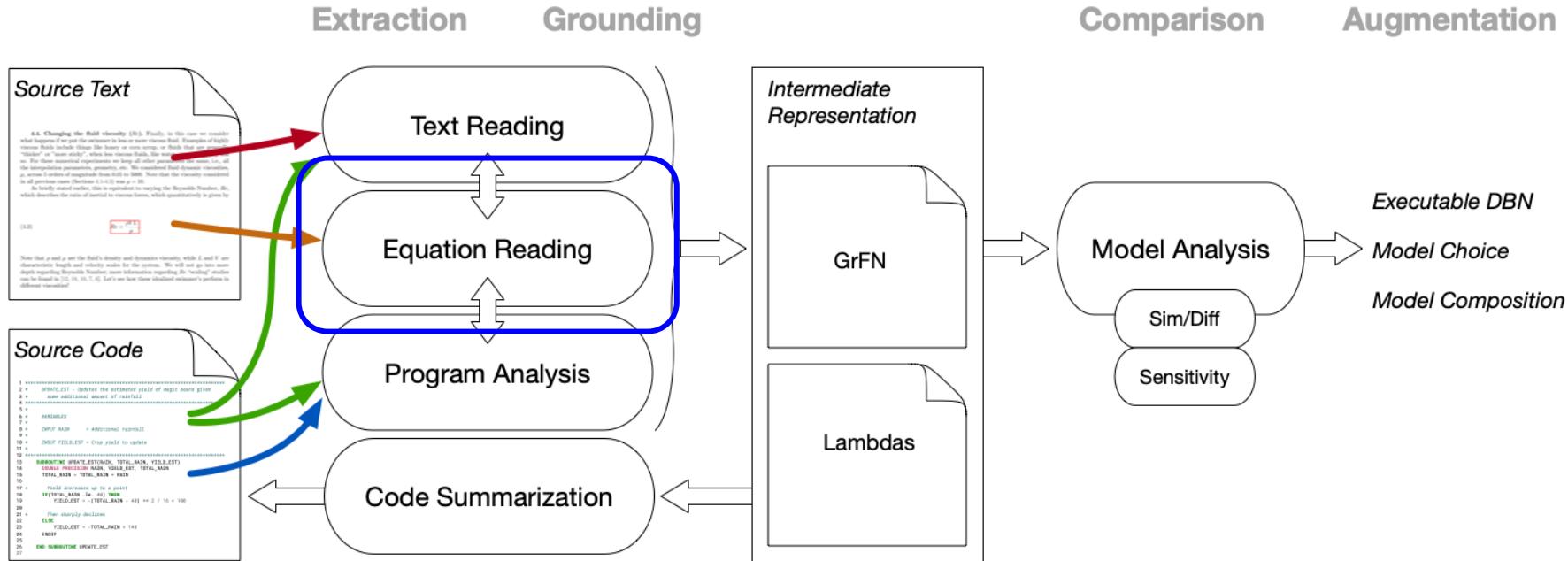


Text Reading Pipeline

Program Analysis



AutoMATES Architecture



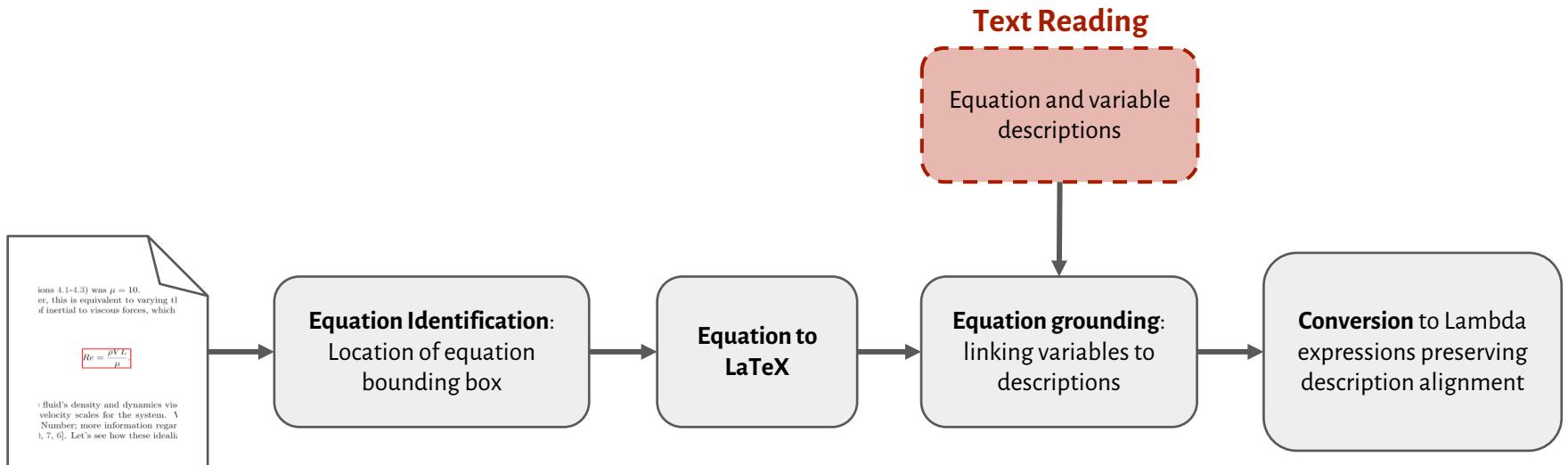
Equation Reading

As briefly stated earlier, this is equivalent to varying the Reynolds Number, Re , which describes the ratio of inertial to viscous forces, which quantitatively is given by

$$(4.2) \quad Re = \frac{\rho V L}{\mu}.$$

Note that ρ and μ are the fluid's density and dynamics viscosity, while L and V are characteristic length and velocity scales for the system. We will not go into more depth regarding Reynolds Number; more information regarding Re “scaling” studies can be found in [12, 18, 10, 7, 6]. Let’s see how these idealized swimmer’s perform in different viscosities!

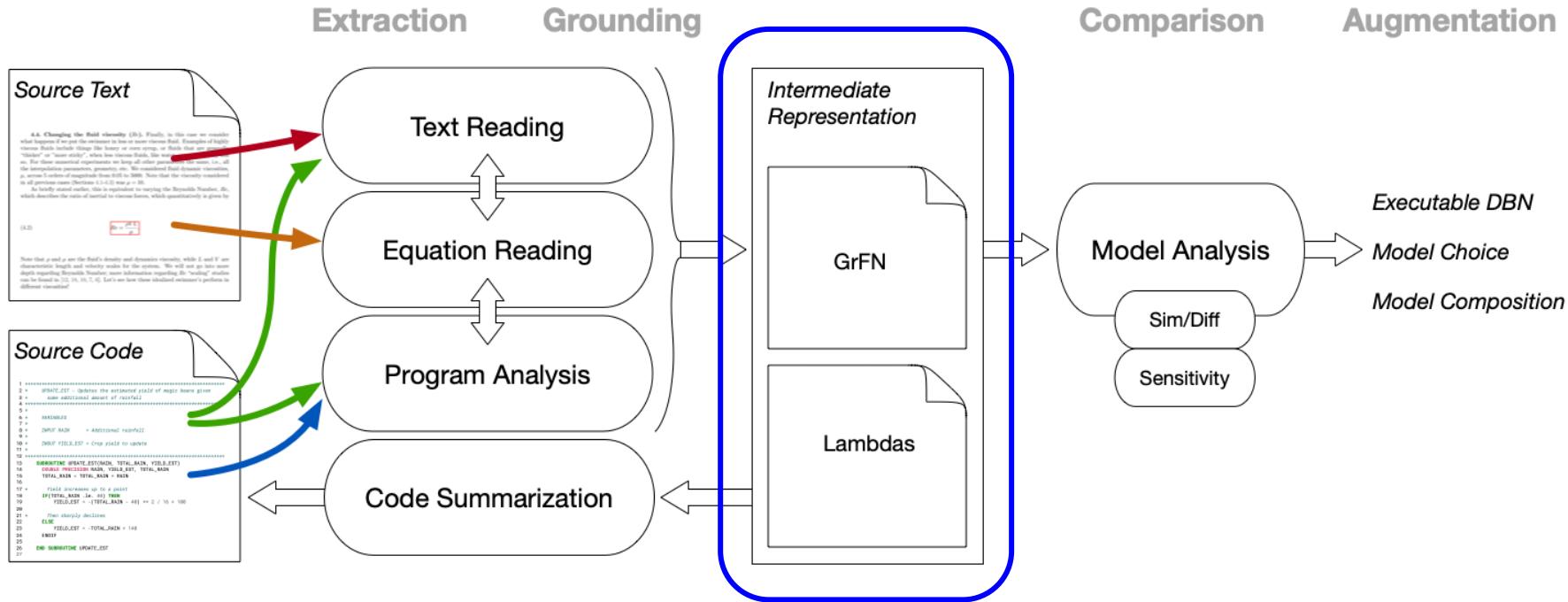
Equation Reading Pipeline



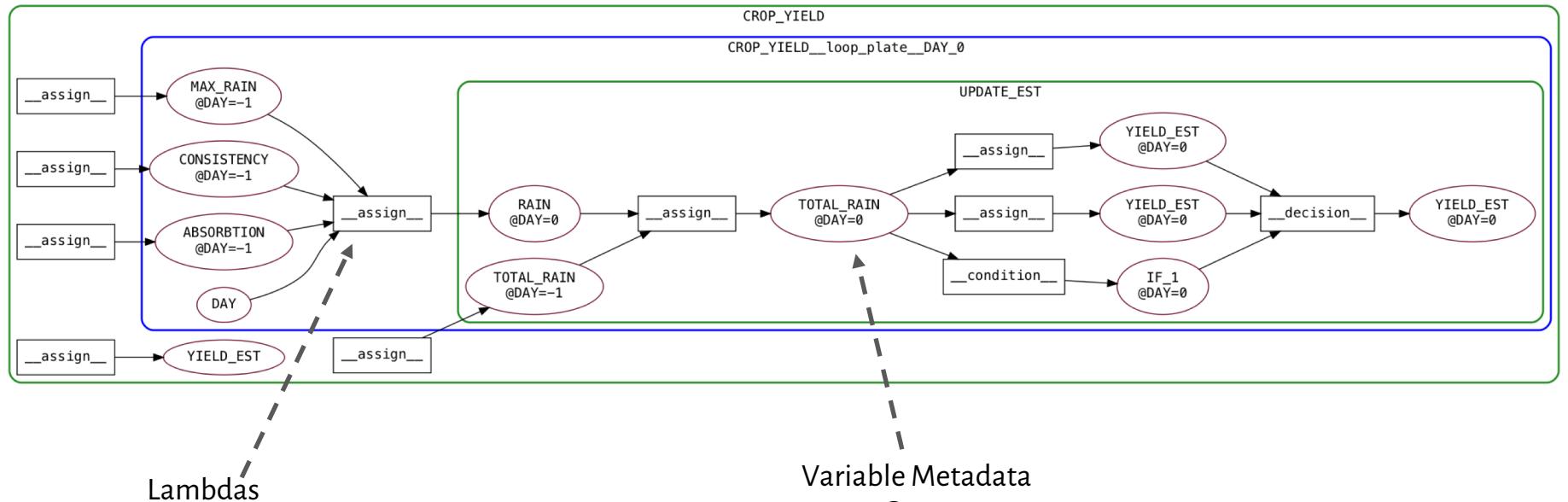
$Re = \frac{\rho V L}{\mu}$

```
def Re(rho, V, L, mu):  
    return (rho * V * L) / mu
```

AutoMATES Architecture

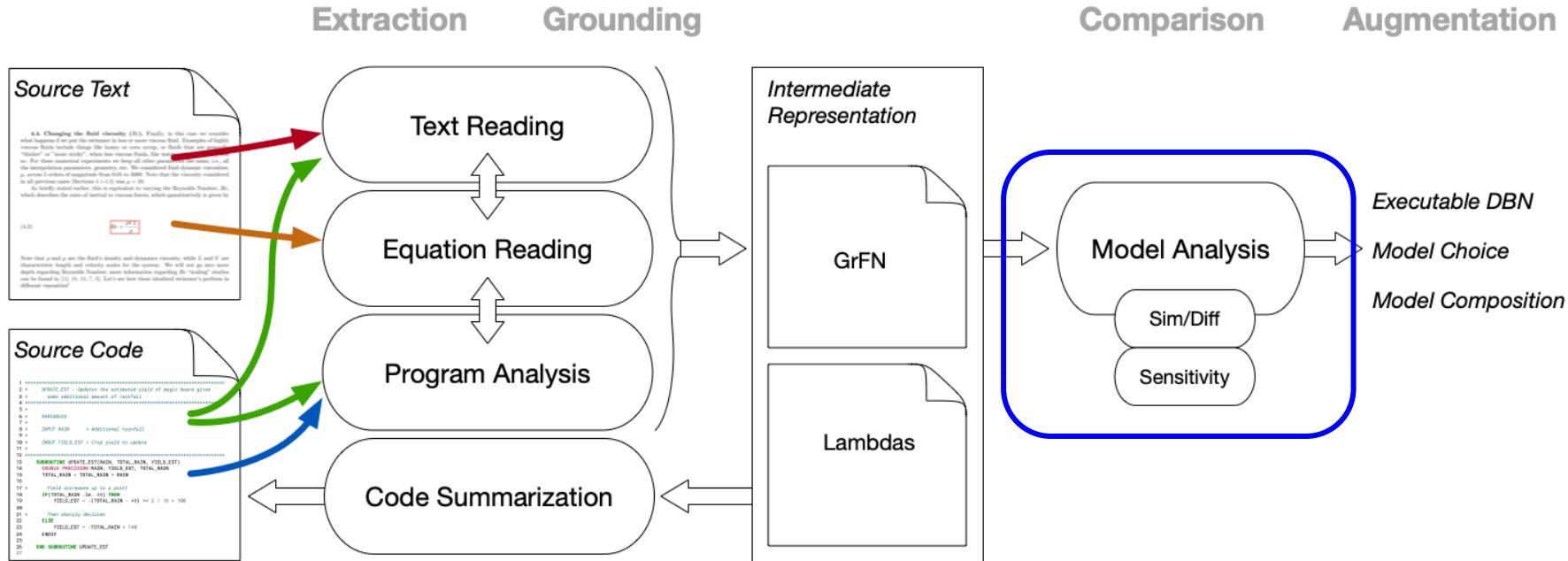


GrFN: Grounded Function Network

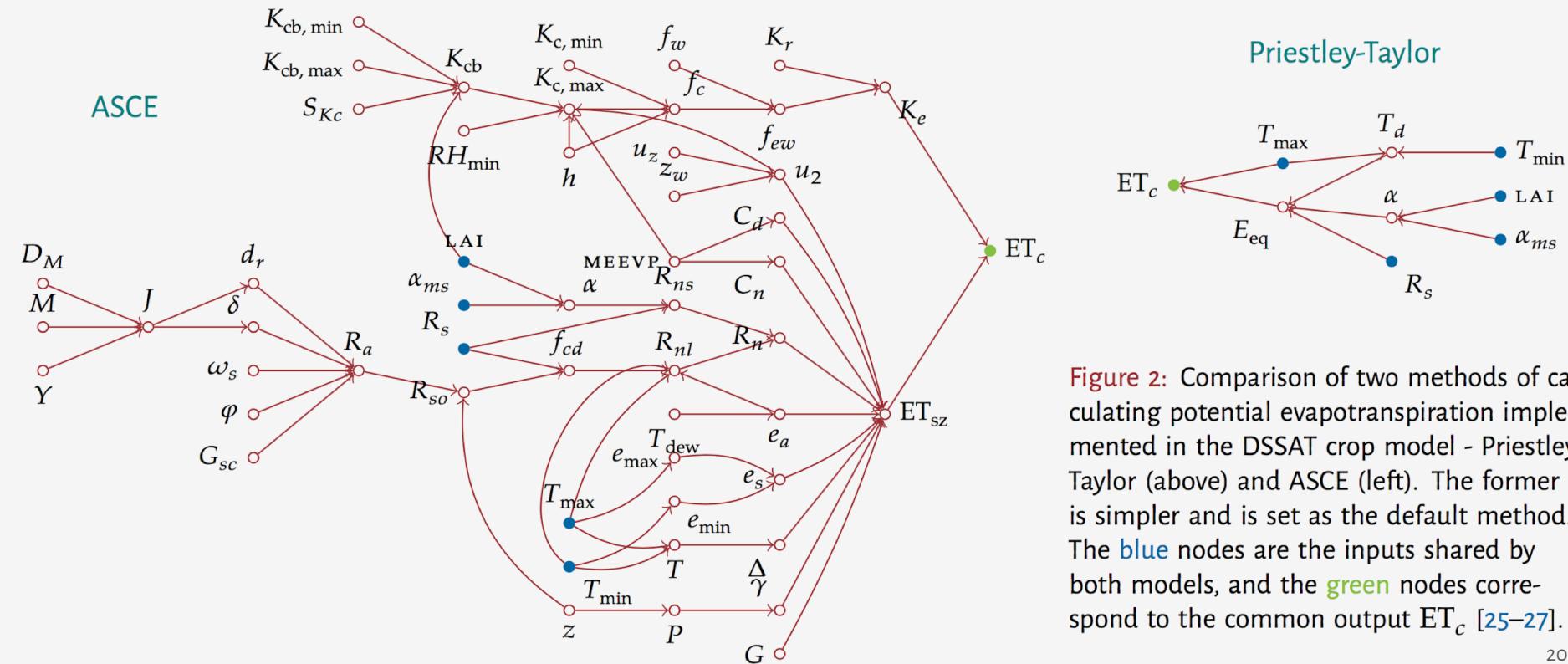


```
def assign_RAIN(MAX_RAIN, CONSISITENCY, ABSORPTON, DAY)
    return (-(DAY - 16) ** 2 / CONSISITENCY + MAX_RAIN)
        * ABSORPTION
```

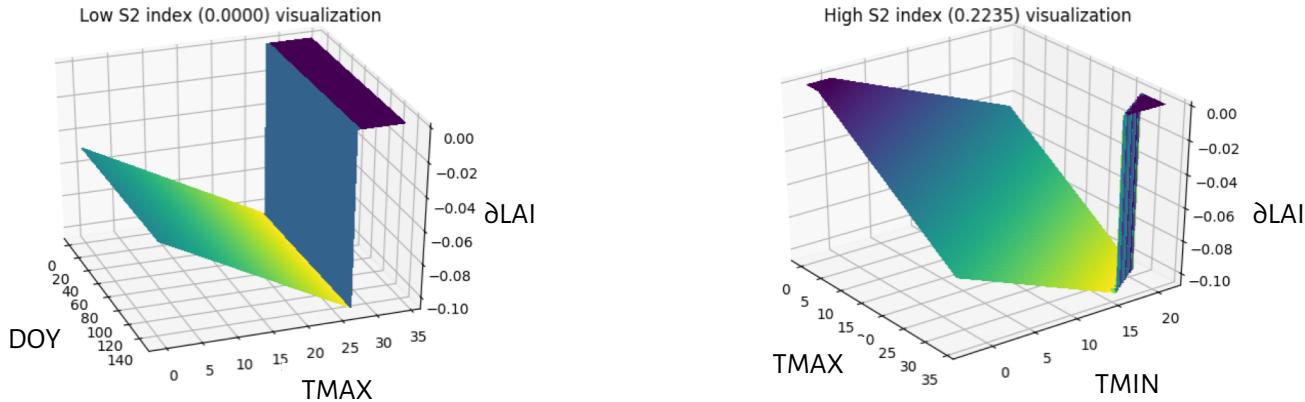
AutoMATES Architecture



Model Analysis: Structural Comparison

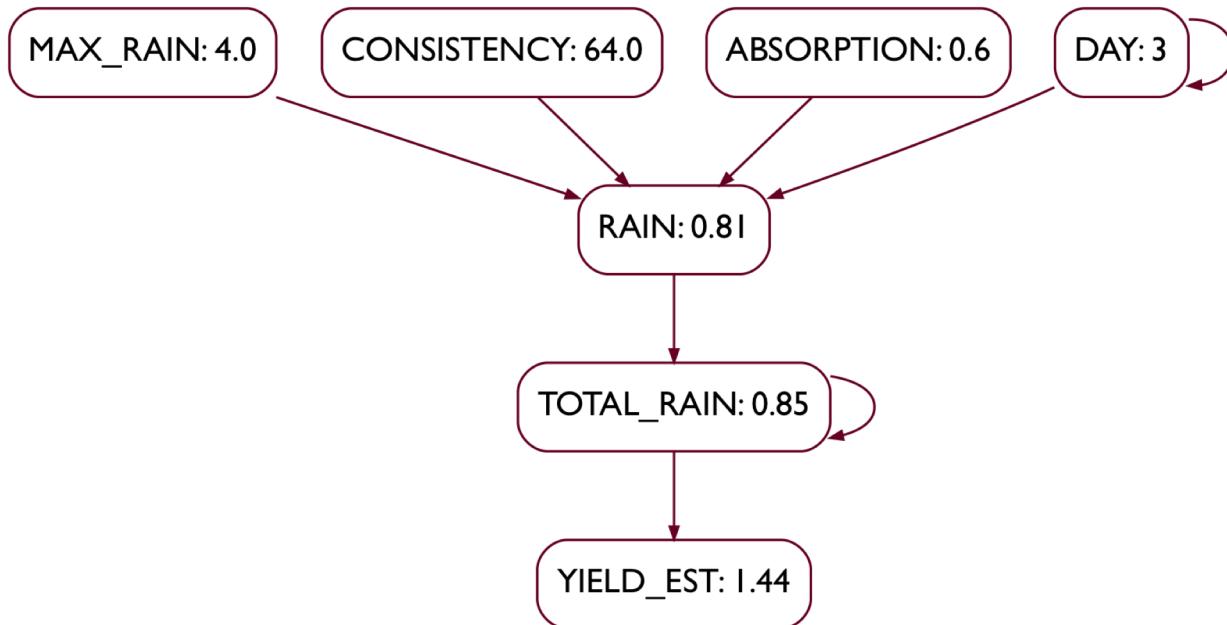


Model Analysis: Sobol Sensitivity Analysis

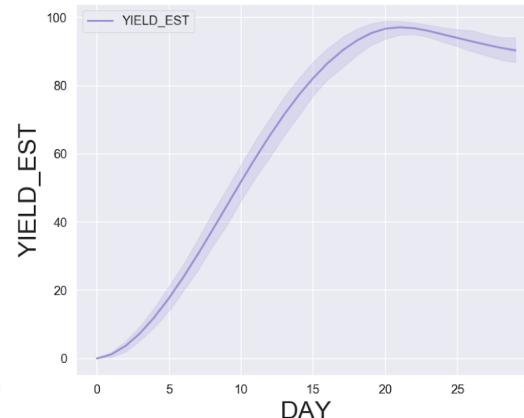
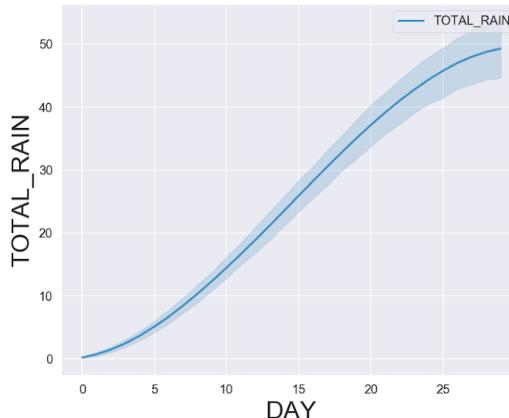
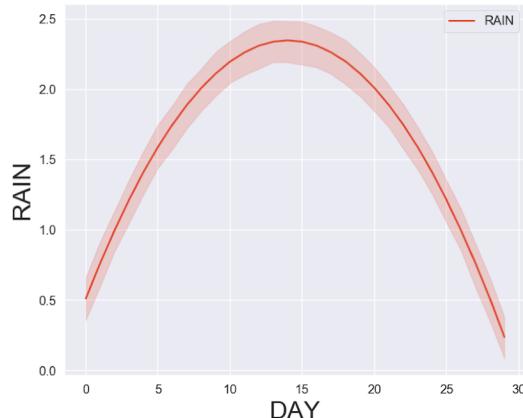
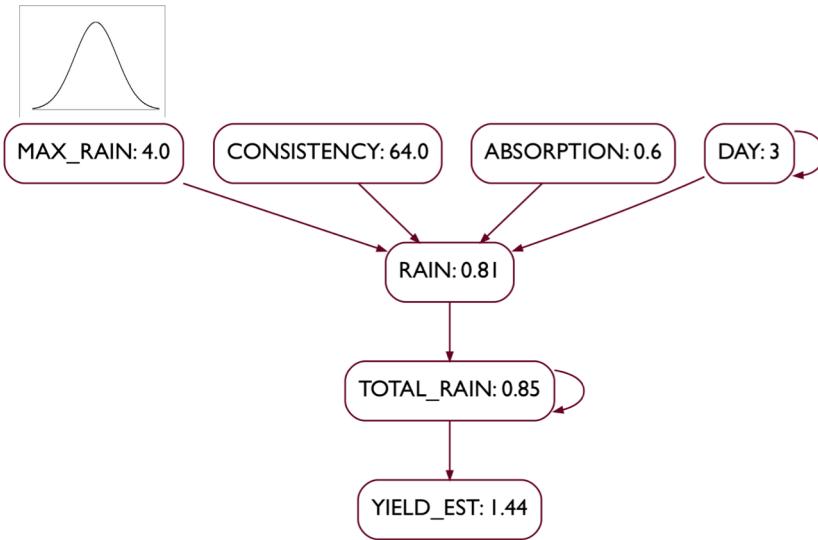


- Sensitivity index computation is done using Sobol sensitivity analysis
- Computing the sensitivity indices allows us to see which variables and which pairs of variables account for the most variance in overall model output

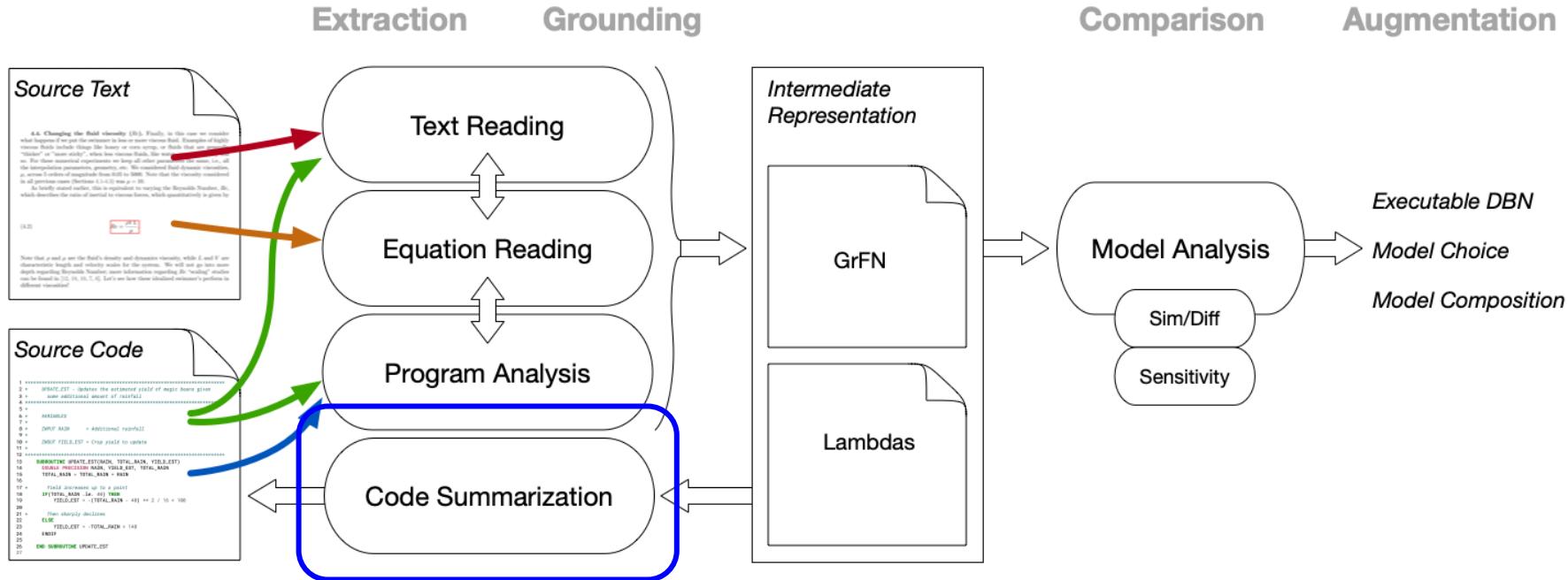
Model Analysis: GrFN as Dynamic Bayesian Network



Model Analysis: DBN



AutoMATES Architecture



Code Summarization: Training Corpus

- Neural network architecture that can encode functions and generate natural language summaries
- Trained using existing, well-documented open source code!



... and more...

Summary

Integration of model semantics
from *text, equations* and *software*
into a *uniform framework* for analysis.

- Automating linking of software to text discourse context
- Model comparison and sensitivity analysis in a uniform framework
- Contextual debugging and model communication
- Facilitate import of source code to libraries (e.g., MINT model store)
- Comparison of **natural language -derived** Causal Analysis Graphs (CAG^{NL} s) to **software** CAGs (CAG^S s)
 - Fill in causal details missing in CAG^{NL} s
 - Expose assumed common-sense knowledge underlying CAG^{NL} s

Thank You!

- Masha Alexeeva
- Pratik Bhandari
- Saumya Debray (co-PI)
- Paul Hein
- Jennifer Kadowaki
- Clay Morrison (PI)
- Adarsh Pyarelal (co-PI)
- Becky Sharp (co-PI)
- Marco Valenzuela-Escárcega (co-PI)

