

# Automated Prior Elicitation for Bayesian Metabolomics Analysis

JSM 2025 | Flexible Prior Elicitation for Bayesian Analysis

Chiraag Gohel

The Rahnavard Lab, The George Washington University

2025-08-06

Automated  
Prior Elicitation  
for Bayesian  
Metabolomics  
Analysis

Chiraag Gohel

Introduction

Methods

Results

# Introduction

# What is metabolomics?

Automated  
Prior Elicitation  
for Bayesian  
Metabolomics  
Analysis

Chiraag Gohel

Introduction

Methods

Results

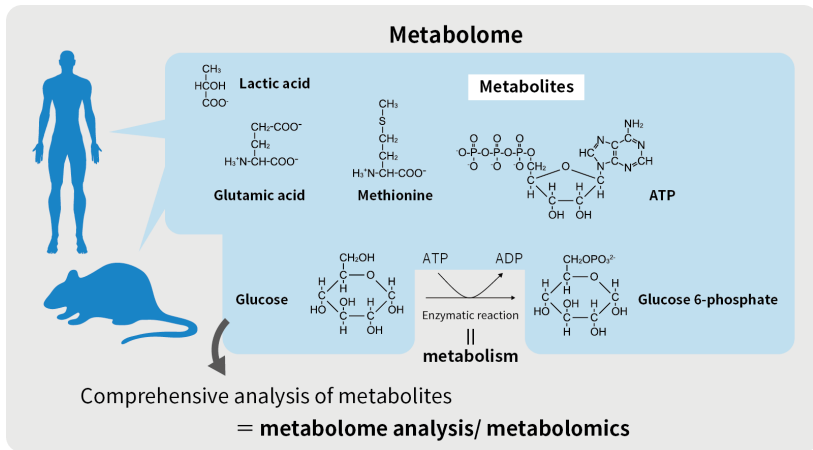


Figure 1: From Human Metabolome Technologies

# Effect size drives biological insight

Automated  
Prior Elicitation  
for Bayesian  
Metabolomics  
Analysis

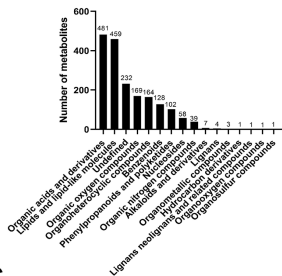
Chiraag Gohel

Introduction

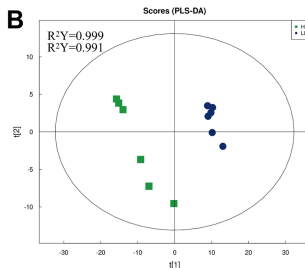
Methods

Results

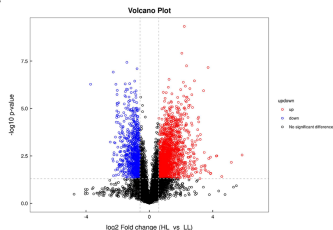
A



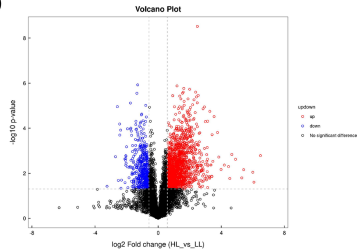
B



C



D



# Traditional testing lacks power

Automated  
Prior Elicitation  
for Bayesian  
Metabolomics  
Analysis

Chiraag Gohel

Introduction

Methods

Results

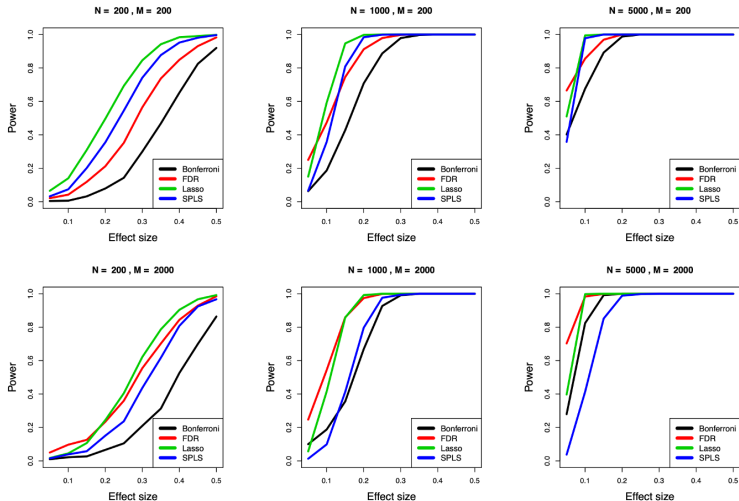


Figure 2: Honglin M. et al. (2022). "Quantitative Comparison of Statistical Methods for

# Prior work

Automated  
Prior Elicitation  
for Bayesian  
Metabolomics  
Analysis

Chiraag Gohel

Introduction

Methods

Results

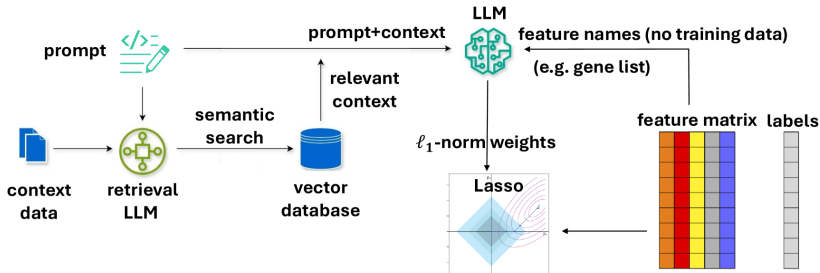
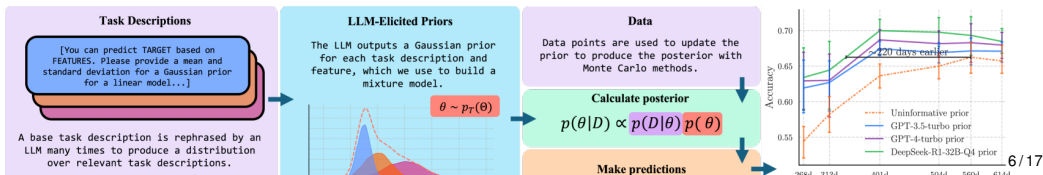


Figure 3: Zhang, E. et al. (2025), “LLM-Lasso: A Robust Framework for Domain-Informed Feature Selection and Regularization,” arXiv.



Automated  
Prior Elicitation  
for Bayesian  
Metabolomics  
Analysis

Chiraag Gohel

Introduction

**Methods**

Results

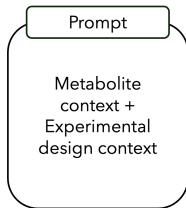
# Methods

# Prior elicitation framework overview



**Human Metabolome Database:** Online database of small molecule metabolites found in the human body

Get  
metabolite  
information



$$\beta_{glucose} \sim N(1.6, 0.4)$$

**Metabolite Context Enrichment:** HMDB database integration for biological context (pathways, functions, disease associations)

**Multi-Model LLM Support:** OpenAI (GPT-4o, O3), Google (Gemini 2.0/2.5), with caching and fallback mechanisms

**Qualitative-to-Numerical Mapping:** Conservative/moderate strength mappings with magnitude-driven effect sizes and confidence-calibrated uncertainties

**Hierarchical Bayesian Modeling:** LLM informed metabolite grouping with



# LLM prior elicitation process

## Step 1: LLM analyzes metabolite + study context

$$\text{LLM}(\text{metabolite, condition}) \quad \{d_j, m_j, c_j, r_j\}$$

where  $d_j$  {increase, decrease, unchanged} is predicted direction,  $m_j$  {small, moderate, large} is predicted magnitude,  $c_j$  (0, 1) is confidence level, and  $r_j$  is a string representing the rationale.

## Step 2: Map qualitative predictions to numerical priors

$$\mu_j^{\text{LLM}} = f(m_j, d_j)$$

$$\sigma_j^{\text{LLM}} = f(c_j)$$

## Step 3: Use as informative priors in Bayesian model $j \sim N(\mu_j^{\text{LLM}}, \sigma_j^{\text{LLM}})$

# Priors

**LLM Priors:**  $j \sim N(\mu_j^{\text{LLM}}, \sigma_j^{\text{LLM}})$

where  $\mu_j^{\text{LLM}}$  and  $\sigma_j^{\text{LLM}}$  are derived from LLM predictions:

$$\mu_j^{\text{LLM}} = m_j \cdot \text{sign}(d_j) \quad (1)$$

$$\sigma_j^{\text{LLM}} = f(c_j) \quad (2)$$

## Conservative Mapping

$m_j \in \{0.08, 0.15, 0.25\}$  for {small, moderate, large}

$f(c_j) \in \{0.5, 0.7, 0.9\}$  for {high, med, low} confidence

**Moderate Mapping**  $m_j \in \{0.12, 0.22, 0.35\}$  for {small, moderate, large}

$f(c_j) \in \{0.3, 0.5, 0.7\}$  for {high, med, low} confidence

# LLM-Informed Hierarchical Prior

Automated  
Prior Elicitation  
for Bayesian  
Metabolomics  
Analysis

Chiraag Gohel

Introduction

Methods

Results

Group metabolites by LLM predictions and use intelligent pooling:

$$\begin{aligned} \text{Group means}_g & \sim N(\mu_g^{\text{LLM}}, 3.0) \\ j & \sim N(\text{Group means}_{g[j]}, 2.0) \end{aligned}$$

where group  $g$  is mapped to  $\mu_g^{\text{LLM}}$  as follows:

$$\mu_g^{\text{LLM}} = \begin{cases} 0.1, & \text{if } g = \text{decrease,} \\ 0.0, & \text{if } g = \text{unchanged,} \\ +0.1, & \text{if } g = \text{increase.} \end{cases}$$

# Modeling

Automated  
Prior Elicitation  
for Bayesian  
Metabolomics  
Analysis

Chiraag Gohel

Introduction

Methods

Results

All Bayesian models use the same log-link GLM structure with different prior specifications:

$$y_{ij} \sim N(\mu_{ij}, \frac{2}{j}) \quad (3)$$

$$\log(\mu_{ij}) = \mu_j + \beta_j x_i \quad (4)$$

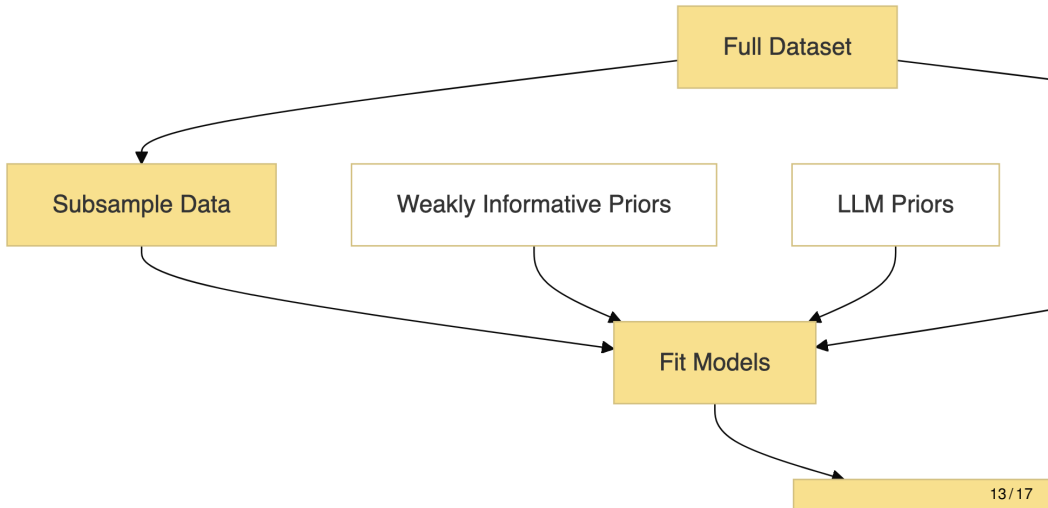
$$\mu_j \sim N(\log(y_j + c), 1.0) \quad (5)$$

$$\beta_j \sim \text{HalfNormal}(0.5) \quad (6)$$

where  $y_{ij}$  is abundance for sample  $i$  and metabolite  $j$ ,  $x_i \in \{0, 1\}$  is group indicator,  $\mu_j$  represents the natural log fold change (lnFC) for metabolite  $j$ , and  $c$  is a small constant to avoid  $\log(0)$ .

# Simulation Study

## Empirical Monte-Carlo Subsampling



Automated  
Prior Elicitation  
for Bayesian  
Metabolomics  
Analysis

Chiraag Gohel

Introduction

Methods

Results

Automated  
Prior Elicitation  
for Bayesian  
Metabolomics  
Analysis

Chiraag Gohel

Introduction

Methods

Results

# Results

# LLM-informed priors improve recovery

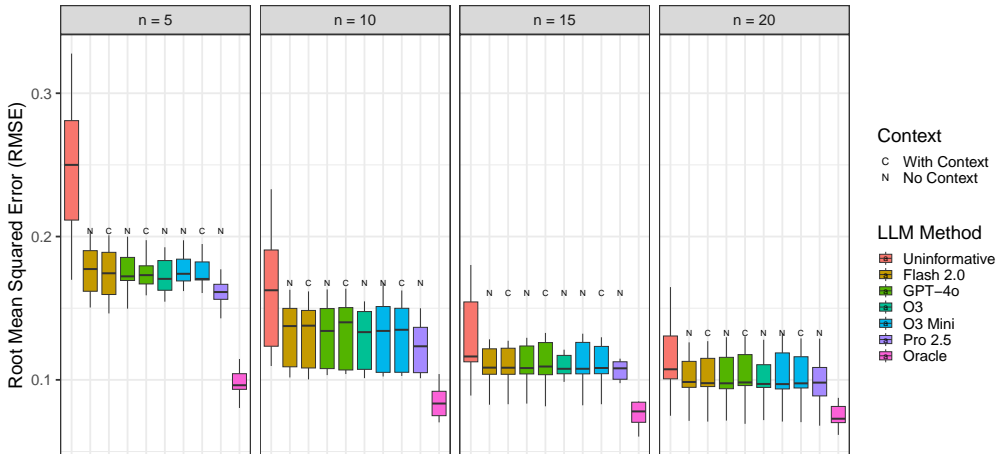
Automated  
Prior Elicitation  
for Bayesian  
Metabolomics  
Analysis

Chiraag Gohel

Introduction

Methods

Results



# LLM Informed estimators are finite-sample efficient

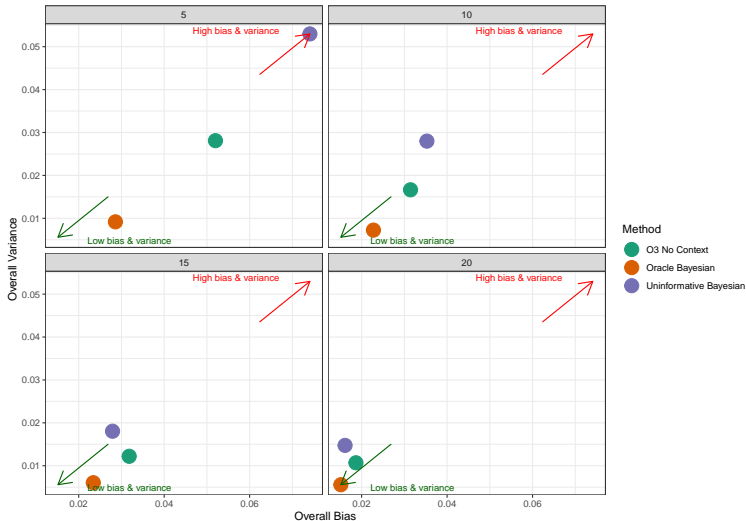
Automated  
Prior Elicitation  
for Bayesian  
Metabolomics  
Analysis

Chiraag Gohel

Introduction

Methods

Results





# Summary

Automated  
Prior Elicitation  
for Bayesian  
Metabolomics  
Analysis

Chiraag Gohel

Introduction

Methods

Results

**LLM Prior Elicitation Works:** Automated biological knowledge extraction via LLMs produces informative priors for Bayesian metabolomics analysis.

**Mapping Strategy Matters:** Magnitude-driven effect sizes and confidence-calibrated uncertainties are crucial for translating qualitative LLM insights into effective numerical priors.

**Added Context May Not Matter:** Including biological context from the HMDB in LLM prompts did not significantly improve prior performance in this study.

**Performance is Model Agnostic:** Different LLMs (OpenAI, Google) yielded similar results, indicating robustness across models.

**Practical Impact:** Method particularly valuable for small sample studies ( $n=5-20$ ) where traditional statistical approaches struggle with high-dimensional metabolomics data.

**Future Directions:** Integration of other databases, alongside more sophisticated mapping approaches and historical data.