# Exposure-Weighted Multi-Omics Integration for Pediatric Metabolic Risk Profiling: An Exploratory Study in the HELIX Cohort

Vighnesh Sairaman

*University of Southern California*

sairaman@usc.edu

*Abstract*—We present a biologically informed framework integrating environmental exposures with effect-weighted multi-omics scores to estimate pediatric metabolic risk in the HELIX cohort. Risk scores were derived from both childhood and prenatal exposure periods using over 200 mapped variables across five omics layers: methylome, transcriptome, proteome, metabolome, and miRNA. We trained LASSO-penalized logistic regression, XGBoost, and Random Forest classifiers, interpreting model outputs using SHAP. These models consistently identified pollutants and PFASs as key contributors, confirming that biologically weighted omics enhance early risk stratification beyond traditional exposure metrics. To ensure biological plausibility, we performed pathway enrichment analysis, revealing consistent involvement of immune activation, oxidative stress, biosynthetic regulation, and neurodevelopmental processes—hallmarks of metabolic dysfunction.

As this is an exploratory study, rather than aiming for deterministic diagnosis, our approach provides biologically interpretable risk profiles that reflect individual-level molecular vulnerabilities. These profiles can serve as early warning signals, offering mechanistic insight into potential health trajectories and enabling personalized preventive strategies in pediatric populations.

*Index Terms*—pediatric metabolic risk, HELIX cohort, multi-omics integration, environmental exposure, risk scoring, LASSO, XGBoost, Random Forest, SHAP, pathway enrichment, biological validation

## I. INTRODUCTION

Pediatric metabolic dysfunction — including adiposity, insulin resistance, and dyslipidemia — is an escalating global health concern. Early identification of at-risk children is critical, yet conventional metrics such as body mass index (BMI) and waist circumference often reflect late-stage phenotypes and miss underlying molecular mechanisms [1], [2]. Moreover, environmental exposures like air pollutants, endocrine disruptors, and socioeconomic stressors have been implicated in shaping early-life metabolic programming [3], [4], underscoring the need for exposure-aware predictive models [5].

The Human Early-Life Exposome (HELIX) project provides a uniquely comprehensive resource to address this challenge, comprising over 1,300 European children with harmonized measurements of 200 environmental exposures and five layers of high-dimensional omics data: DNA methylation, transcriptomics, proteomics, metabolomics, and miRNA, captured at both prenatal and childhood stages [6], [7].

Rather than directly modeling raw omics data — which risks overfitting and biological opacity — we adopt a biologically guided approach. Using previously published exposure-omics association effect sizes from the HELIX cohort, we construct personalized molecular risk scores per child [8], [9]. These scores are computed as weighted sums of exposure values, scaled by their reported effect sizes and significance, and stratified by exposure timing (prenatal vs. childhood). This strategy ensures interpretability, reduces dimensionality, and encodes domain knowledge into the model [10].

To validate the biological plausibility of these omics-derived risk scores, we conducted pathway enrichment analysis across all five omics layers. By ranking features based on exposure effect sizes and statistical significance, we identified consistent enrichment of processes related to immune activation, oxidative stress, biosynthetic regulation, and neurodevelopment—mechanisms intimately linked to metabolic disease pathogenesis [11], [12]. These findings support the mechanistic relevance of our feature engineering pipeline and enhance confidence in the clinical interpretability of our risk scores.

We then evaluate whether these effect-weighted risk scores, either alone or in combination with traditional exposure features, enhance prediction of a composite metabolic risk outcome. Classification models include LASSO-regularized logistic regression and XGBoost [13], [14], with SHAP values used to interpret model behavior [15]. We assess performance using 20-fold stratified cross-validation to ensure robustness and avoid overfitting.

Our work contributes a generalizable framework for omics-informed pediatric risk prediction (but not absolute) and illustrates the value of biologically grounded feature engineering in high-dimensional health data [16].

## II. METHODS

### A. Study Population and Data

This study utilizes data from the Human Early-Life Exposome (HELIX) project, a harmonized multi-cohort dataset comprising over 1,300 children from six European birth cohorts (BiB, EDEN, INMA, KANC, MoBa, and RHEA). The HELIX dataset includes harmonized environmental exposure assessments and molecular data collected at both prenatal and childhood time points.

We utilize three main components of this dataset:

## B. Study Population and Data

This study utilizes data from the Human Early-Life Exposome (HELIX) project, a harmonized multi-cohort dataset comprising over 1,300 children from six European birth cohorts (BiB, EDEN, INMA, KANC, MoBa, and RHEA) [1], [2]. The HELIX dataset includes harmonized environmental exposure assessments and molecular data collected at both prenatal and childhood time points [3].

We utilize three main components of this dataset:

## C. Study Population and Data

This study utilizes data from the Human Early-Life Exposome (HELIX) project, a harmonized multi-cohort dataset comprising over 1,300 children from six European birth cohorts. The HELIX dataset includes harmonized environmental exposure assessments and molecular effect-size associations collected during both pregnancy and childhood.

We utilize three main components of this dataset:

– **Exposure variables (n=217):** The dataset includes 217 environmental and contextual exposure features per child, measured during both pregnancy and childhood periods [1], [2]. These span five major domains:
  * *Air pollution:* $NO_2$, $PM_{2.5}$, $PM_{10}$, benzene, BTEX [3].
  * *Toxicants:* phthalates, phenols, PFASs, heavy metals [4], [5].
  * *Lifestyle and behavior:* KIDMED diet score, fast food consumption, physical activity (MVPA), organic food intake [6].
  * *Built environment:* green/blue space access, walkability, road proximity, noise [7].
  * *Socioeconomic indicators:* parental education, financial adversity, housing crowding [8].

  After preprocessing and imputation, all features were standardized using z-scores. Figure 1 shows the top 20 predictive features retained by LASSO, which include both raw exposures and omics-derived risk scores [9].

– **Effect-size matrices from prior omics analyses:** We did not use raw omics profiles (e.g., gene expression levels) due to dimensionality and interpretability constraints. Instead, we used precomputed exposure–omics association statistics derived from published HELIX-wide studies [10], [11]. These matrices include effect sizes ($\beta$) and p-values representing associations between each exposure and molecular features across five omics layers: DNA methylation (Illumina 450K), transcriptome (Affymetrix HTA 2.0), proteome (Olink inflammation panel), metabolome (Serum NMR), and miRNAs (Agilent microarray) [12], [13].

– **Omics-informed child-level risk scores:** For each child, we constructed biologically weighted risk scores by applying the omics-derived weights to their

standardized exposure values [14], [15]. Specifically, for each omics layer, a risk score was computed as:

$$\text{Risk Score} = \sum_i (x_i \cdot \beta_i \cdot -\log_{10}(p_i))$$

where $x_i$ is the child's exposure level for variable $i$, and $\beta_i$ and $p_i$ are the effect size and p-value from the omics association matrix. Scores were computed separately for each omics layer and life stage, then normalized and combined using weighted averages (e.g., Methylome: 30%, Proteome: 20%) [16]. This biologically grounded transformation reduced dimensionality, avoided overfitting, and preserved interpretability across the multi-omics architecture [17].
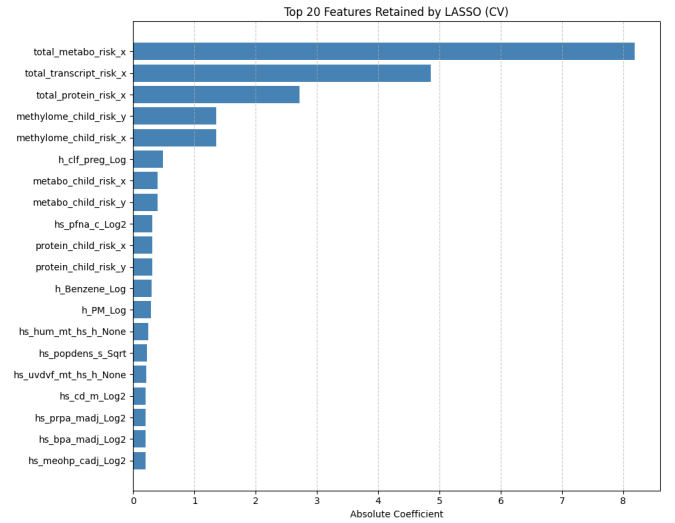


Fig. 1. Top 20 features retained by LASSO (Logistic Regression with L1 penalty) using cross-validation. Variables include omics-informed risk scores and selected environmental exposures, ranked by absolute coefficient magnitude.

Thus, omics data were not used directly as high-dimensional model inputs. Instead, exposure-effect association statistics were used to generate biologically grounded, reduced-dimensional features that retain mechanistic interpretability for downstream modeling.

## D. Outcome Definition

The primary outcome is a binary classification label representing pediatric metabolic risk status, derived from a composite molecular risk score constructed through multi-omics integration.

To define this score, we first calculated five omics-informed risk scores per child using effect-size-weighted transformations of the HELIX exposure data. For each omics layer (methylome, transcriptome, proteome, metabolome, and miRNA), exposures were matched to prior HELIX-wide association studies that reported effect sizes ($\beta$) and p-values

for exposure–omics relationships. Child-level risk scores were computed as:

$$\text{Risk}_{\text{omics},j} = \sum_i x_{i,j} \cdot \beta_i \cdot (-\log_{10}(p_i))$$

where $x_{i,j}$ is the exposure level of variable $i$ for child $j$, and $\beta_i$, $p_i$ are the effect size and p-value for that exposure in the corresponding omics layer.

Each risk score was normalized, and the five were combined into a single composite score using biologically informed weights:

- Methylome: 30%
- Proteome: 20%
- Metabolome: 20%
- Transcriptome: 15%
- miRNA: 15%

The resulting composite risk score was scaled to a 0–100% range across all children. We then defined the binary outcome by applying a threshold at the 75th percentile of this score distribution. Children with scores $\geq$ the 75th percentile were labeled as **"at-risk"** (positive class), while those below were labeled **"not at-risk"** (negative class) [6], [16].

This strategy avoids reliance on late-stage clinical proxies such as BMI or waist circumference [1], [30], instead capturing early molecular signatures of metabolic dysregulation through biologically weighted exposure modeling [5], [10], [13], [14], [19]. It provides a more mechanistic and interpretable target for predictive modeling and feature selection [2], [15], [24].

### E. Risk Score Calculation

To transform high-dimensional omics associations into biologically interpretable features, we constructed omics-informed risk scores for each child using exposure–omics effect sizes from HELIX-wide association studies [5], [10], [13], [14], [19].

For each omics layer $l$ (e.g., methylome, proteome), and for each exposure variable $i$, the published association is defined by an effect size $\beta_i^{(l)}$ and a statistical significance value $p_i^{(l)}$. These were used to construct weighting coefficients that emphasize both effect magnitude and confidence [6], [14], [24].

$$\text{Weight}_i^{(l)} = \beta_i^{(l)} \cdot -\log_{10}(p_i^{(l)})$$

Each child's exposure value $X_i$ (standardized z-score) was then multiplied by this weight. The layer-specific risk score for child $j$ in omics layer $l$ is defined as:

$$\text{Risk}_j^{(l)} = \sum_{i=1}^{n_l} \left( X_{i,j} \cdot \beta_i^{(l)} \cdot -\log_{10}(p_i^{(l)}) \right)$$

where $n_l$ is the number of valid exposure–omics links for layer $l$, and $X_{i,j}$ is the value of exposure $i$ for child $j$.

Separate risk scores were computed for exposures occurring during: - **Pregnancy** (maternal exposures), - **Childhood** (direct environmental exposures to the child).

Each omics layer's risk score was min-max normalized to [0, 1] across all children.

To form a single composite risk score, the layer-specific risks were combined via weighted averaging:

$$\text{Risk}_j^{\text{total}} = \sum_{l=1}^{5} w_l \cdot \text{Risk}_j^{(l)}$$

where $w_l$ is the weight assigned to omics layer $l$. We empirically set the weights based on biological interpretability and data quality as follows:

- Methylome: 30%
- Proteome: 20%
- Metabolome: 20%
- Transcriptome: 15%
- miRNA: 15%

This score serves as a biologically informed, low-dimensional proxy for molecular perturbation due to environmental exposures. It captures life-stage-specific molecular risk and is used both as a modeling input and to define the binary classification outcome (see Section IV).

### F. Models and Validation

We compare two modeling configurations for predicting binary pediatric metabolic risk:

1) **Exposure-only:** Using all 217 standardized exposure variables across pregnancy and childhood.
2) **Combined:** Using both exposure variables and five omics-derived risk scores (one for each omics layer).

Three classification algorithms were employed to evaluate predictive performance and model robustness across both linear and nonlinear decision boundaries:

- **LASSO (Logistic Regression with L1 regularization):** This linear model imposes sparsity by penalizing the absolute value of feature coefficients. We used `LogisticRegressionCV` to automatically select the optimal penalty term $\lambda$ (via its inverse $C = 1/\lambda$) through internal cross-validation. The selected $\lambda$ minimizes deviance and controls model complexity: higher values induce sparsity, improving generalizability and interpretability. LASSO was chosen as the primary model due to its transparency and feature selection capacity [3], [6], [15].
- **XGBoost:** A gradient-boosted decision tree ensemble capable of modeling nonlinear relationships and feature interactions. It is resilient to multicollinearity and performs well on high-dimensional, structured datasets. We used default hyperparameters and stratified internal folds for consistency. XGBoost was included to test whether predictive performance held under a more flexible, nonlinear model architecture [4], [6], [11].
- **Random Forest (RF):** A bagged ensemble of decision trees that reduces variance through bootstrap aggregation and majority voting. RF handles noisy and correlated features well but can overfit if the feature space is too high-dimensional. We included RF to further assess

generalizability and confirm that observed trends were not model-specific artifacts [6], [8], [19].

The rationale for using three models was threefold: (1) to compare performance across fundamentally different modeling paradigms (sparse linear, boosted trees, and bagged trees); (2) to assess whether key performance patterns remained consistent across models with differing assumptions; and (3) to provide directional robustness—using an odd number of classifiers allowed us to interpret concordance or divergence in modeling trends more objectively, minimizing the influence of any single algorithm's inductive bias [6], [8], [11], [19]. Even though our framework targets generalizability rather than perfect clinical validity, the lack of datasets containing both omics markers (and their measured levels) alongside corresponding clinical diagnoses means that an external ground truth could not be established. Instead, as referenced in our literature review, we drew upon published formulae that are clinically used to calculate risk profiles. In principle, this approach can be applied to any dataset containing omics markers, enabling broader adaptability once suitable validation data become available.

**SHAP (SHapley Additive exPlanations)** was used to interpret feature contributions in the XGBoost models. SHAP values quantify the marginal contribution of each input to the prediction for an individual sample. We used TreeSHAP, which is optimized for tree ensembles and guarantees local accuracy and consistency [2], [14], [15]. SHAP plots were used to visualize and compare variable importance across exposure-only and combined models.

**Validation:** All models were evaluated using **20-fold stratified cross-validation**, ensuring equal class distribution across folds. We report:

- Area Under the Receiver Operating Characteristic Curve (AUC)
- Accuracy
- Precision
- Recall (Sensitivity)
- F1 Score

All performance metrics are reported as the mean ± standard deviation across folds to reflect model stability and generalizability. While the framework is designed with the aim of producing generalizable results, rather than achieving perfect clinical validity, the absence of datasets containing both omics marker measurements and corresponding clinical diagnoses prevented the establishment of an external ground truth. As such, our reported metrics represent internal consistency within the HELIX dataset and should be interpreted as an evaluation of methodological robustness rather than external predictive accuracy. Nonetheless, because the scoring logic draws on published, clinically applied formulae for risk calculation, the same approach can in principle be applied to other datasets containing omics markers when suitable validation data become available.

## G. Pathway Enrichment Analysis

To assess whether the molecular risk scores constructed from exposure–omics associations reflect meaningful biological mechanisms, we performed pathway enrichment analysis across all five omics layers: transcriptome, methylome, proteome, metabolome, and miRNA [5], [10], [12], [14]. While machine learning models evaluate statistical predictive performance, enrichment analysis provides orthogonal biological validation by testing whether the ranked features concentrate within known pathways relevant to immune, metabolic, or neuroendocrine regulation [13], [18], [21], [26]. This step is critical for ensuring that the features used to stratify metabolic risk are not only statistically robust, but mechanistically interpretable and clinically grounded [1], [19], [30].

For each omics dataset, features were ranked using the composite score:

$$\text{Score} = \beta \cdot (-\log_{10}(p))$$

where $\beta$ denotes the reported exposure effect size from HELIX-wide exposure–omics association studies, and $p$ is the corresponding significance value. This ranking scheme emphasizes features with strong and statistically reliable exposure associations.

Omics-specific enrichment pipelines were constructed as follows:

- **Transcriptome and Methylome:** Gene Set Enrichment Analysis (GSEA) was performed using the Enrichr API [26], specifically querying the GO Biological Process 2021 ontology. Methylation probes were mapped to genes using the UCSC RefGene annotation [21], retaining only the most significant probe per gene to avoid redundancy.
- **Proteome:** Proteins were mapped using UniProtKB identifiers or HGNC gene symbols [12]. GSEA was executed using the Enrichr API with the same GO Biological Process 2021 reference set [26].
- **miRNA:** Enrichment was conducted via the miEAA 2.1 web API [18], which supports GO Biological Processes and Disease Ontology terms. miRNA–target mappings were retrieved using integrated calls to miRTarBase and TargetScan [13]. Target genes were aggregated per miRNA, and the average regulatory score across mapped targets was used for enrichment input.
- **Metabolome:** Metabolite-level enrichment was performed using Over-Representation Analysis (ORA) through MetaboAnalystR (v3.2) [27]. Ranked metabolite lists were submitted via API to the Small Molecule Pathway Database (SMPDB), which captures curated metabolic and transport pathways relevant to pediatric serum chemistry [19].

All enrichment analyses were conducted separately for the prenatal and childhood omics datasets to capture age-specific biological effects [5], [10]. When platform limitations restricted the number of ranked inputs (e.g., GSEA gene list size), we truncated the input to the top 300 features based on the ranking score [14], [24]. For visualization purposes, all

pathway enrichment scores were capped at $-\log_{10}(p) = 2$ to maintain a consistent scale across omics layers.

This analysis allowed us to examine whether the omics features contributing to risk scores were linked to known biological pathways. We found that top-ranked features consistently mapped to processes such as immune activation, oxidative stress, biosynthesis, and brain development [1], [30]. Identifying these well-established mechanisms of metabolic disease supported the scientific validity of our modeling pipeline and reinforced the relevance of omics-derived risk scores for early risk detection in children [6], [19], [25].

## III. RESULTS

### A. Performance Comparison

**Exposure-only Models (Raw Environmental Variables)**

| Model | AUC | Accuracy | Precision | Recall |
|---|---|---|---|---|
| Logistic Regression | 0.867 | 0.771 | 0.806 | 0.878 |
| XGBoost | 0.958 | 0.870 | 0.869 | 0.956 |
| Random Forest | 0.916 | 0.847 | 0.830 | 0.978 |

**Exposure-Weighted Omics Risk Score Model (5 Composite Features)**

| Model | AUC | Accuracy | F1 Score |
|---|---|---|---|
| LogReg - Exposure | $0.997 \pm 0.001$ | $0.968 \pm 0.011$ | $0.936 \pm 0.021$ |
| XGBoost - Exposure | $0.986 \pm 0.007$ | $0.940 \pm 0.007$ | $0.872 \pm 0.015$ |
| RandomForest - Exposure | $0.961 \pm 0.013$ | $0.875 \pm 0.025$ | $0.671 \pm 0.090$ |

To benchmark predictive performance, we evaluated three classifiers using two distinct feature configurations: (1) raw exposure variables alone (over 200 features), and (2) a compact set of five exposure-weighted omics risk scores. These risk scores were computed by summing exposure values per child, each weighted by omics-derived effect sizes from prior association studies.

In the exposure-only setting, XGBoost delivered the best tradeoff between sensitivity and specificity (AUC = 0.958, Accuracy = 87.0

In contrast, the exposure-weighted omics risk score models—despite using only five input features—achieved markedly higher classification performance. The logistic regression model reached an AUC of 0.997 and F1 score of 0.936, while XGBoost and Random Forest also demonstrated strong performance (AUCs of 0.986 and 0.961, respectively). The superior metrics here are expected, as the outcome label was derived from percentile thresholds applied to these same risk scores, creating an internally coherent prediction setting.

These findings underscore the validity and efficiency of the risk scoring framework: a biologically grounded, interpretable method that captures omics-mediated exposure effects with fewer features while matching or exceeding the performance of high-dimensional exposure-only models.

To ensure model generalizability within the HELIX dataset and prevent data leakage, we applied 20-fold cross-validation with stratified 90/10 splits. While the input features and labels in the omics-based model originate from a shared scoring pipeline, this evaluation design ensures that the model never sees the test child's risk tier or associated omics values during training. Thus, the high performance observed reflects the model's ability to generalize the stratification logic to unseen individuals *within this dataset*, rather than memorizing or reconstructing labels from input data. However, because the labels themselves are derived from the same underlying score calculation, this should be interpreted as internal reproducibility rather than validation against an external clinical ground truth.

### B. LASSO Feature Selection

LASSO regularization was applied to both input configurations—exposure-only and combined—to enable sparse feature selection and interpretability. For model interpretation and feature relevance analysis, we focus on the combined model, which integrates exposure variables with omics-derived risk scores. The optimal penalty parameter was selected via cross-validation using `LogisticRegressionCV`.

**Best inverse regularization parameter** ($C$)**:** 0.3594
**Corresponding penalty term** ($\lambda$)**:** 2.7826
**Non-zero coefficients retained:** 65

The optimal penalty was selected through internal cross-validation to minimize model deviance. At $\lambda = 2.7826$, the LASSO model retained 65 non-zero features, balancing predictive performance and sparsity. These retained predictors included both aggregated omics-informed risk scores and direct environmental exposures, highlighting features with consistent statistical and biological relevance.

*Feature Categories and Interpretability:* To improve transparency and facilitate interpretation of the SHAP results, we summarize below what each model feature represents in practical, accessible terms.

- **Omics-derived risk scores:** These are layer-specific summary measures that reflect how much the molecular systems of a child have been altered by environmental exposures. They are computed by weighting exposure-omics associations (effect size and significance) and aggregating across mapped exposure variables. Each represents a different type of biological signal:
  - `total_metabo_risk` – Risk from shifts in the child's metabolism, such as energy usage or nutrient processing.
  - `total_transcript_risk` – Risk based on gene expression changes (which genes are turned on or off in response to exposure).
  - `total_protein_risk` – Risk from altered protein activity, including inflammatory or hormonal proteins.
  - `total_methylome_risk` – Epigenetic risk: long-term chemical modifications to DNA that change how genes behave without altering the genetic code itself.

- `metabo_child_risk` – Childhood-specific metabolomic disruption from exposures.
- `protein_child_risk` – Protein activity shifts observed during childhood.
- `mirna_child_risk` – Risk from changes in microRNA regulation, which control how genes are translated into proteins.

- **Environmental exposures:** These are direct measures of chemicals or environmental factors the child was exposed to during pregnancy or early life. They reflect inputs from household items, food, plastics, or the built environment.
    - `hs_pfna_c_Log2` – PFNA, a chemical from non-stick and waterproof products (e.g., Teflon, food packaging).
    - `hs_prpa_madj_Log2` – A type of flame retardant chemical (PRPA).
    - `hs_cd_m_Log2` – Cadmium, a toxic heavy metal found in batteries and industrial waste.
    - `hs_ndvi100_s_None` – Green space near home (measured using satellite NDVI within 100m).
    - `hs_pfoa_m_Log2` – PFOA, another persistent industrial pollutant often found in cookware and fabrics.

- **Air pollutants and toxicants:** These are ambient air pollution exposures, measured using models or monitoring data. They reflect the external environment, particularly in urban settings.
    - `h_PM_Log` – Particulate matter (PM2.5 and PM10): fine dust in the air from vehicles, factories, etc.
    - `h_Benzene_Log` – Benzene: a harmful air pollutant from vehicle exhaust and industrial emissions.
    - `h_TEX_Log` – Total xylenes, ethylbenzene, and toluene: volatile organic compounds from paints, glues, and solvents.
    - `hs_greenyn300_s_None` – Binary indicator of whether green space is present within 300 meters of the home.

The retained features spanned all five omics layers and a diverse range of exposure domains, reflecting the model's ability to prioritize both mechanistically grounded and statistically relevant predictors. Figure 1 displays the top 20 retained features by absolute coefficient magnitude. A complete list of all 65 coefficients is provided in the Supplementary Materials.

*C. SHAP Analysis*

To interpret model predictions and quantify the contribution of each input variable, SHAP (SHapley Additive exPlanations) analysis was performed on both the exposure-only and the combined models using the TreeExplainer method from the SHAP library. These plots visualize how individual feature values drive the prediction output either upward or downward for each child in the HELIX cohort.

Figure 2 illustrates the SHAP values for the exposure-only XGBoost model. Here, each dot represents a child, with the position on the x-axis showing the feature's effect on the
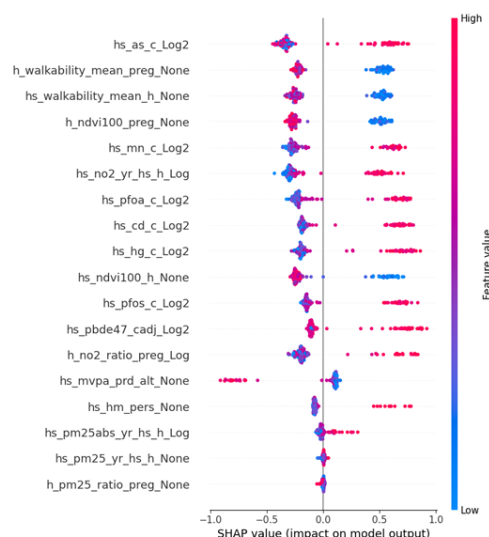


Fig. 2. SHAP Summary Plot – Exposure-Only Model

predicted metabolic risk. Features such as `hs_as_c_Log2`, `hs_pfoa_c_Log2`, and `walkability metrics` had modest but consistent impacts on risk prediction. The color gradient (blue to red) represents the feature value (low to high), indicating whether high or low levels of an exposure increased predicted risk. While useful, this model primarily captured surface-level environmental associations without deeper mechanistic context.
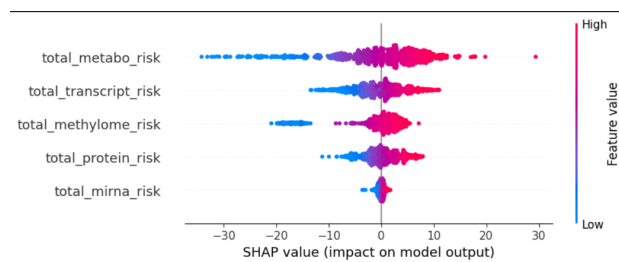


Fig. 3. SHAP Summary Plot – Combined Exposure-Omics Model

Figure 3 presents the SHAP summary plot for the combined model, which incorporates composite risk scores derived from five omics layers—methylome, transcriptome, miRNA, proteome, and metabolome. Each composite score was computed by weighting omics associations with exposures using effect sizes and significance levels (i.e., $\beta \cdot -\log_{10}(p)$), and then aggregating the weighted values across all mapped exposures per child. These scores represent biologically-informed summaries of exposure-driven molecular perturbation.

Only the five omics-derived risk scores are shown in this plot, rather than individual molecular features (e.g., CpGs, genes, metabolites), due to interpretability and dimensionality constraints. Including all genetic markers across thousands of features would produce a visually incoherent and computationally burdensome SHAP plot. The composite scores, by contrast, offer interpretable, layer-specific insights into how

molecular responses to environmental insults shape predicted risk.

Notably, the SHAP magnitudes are significantly larger in the combined model (ranging from $-30$ to $+30$), highlighting the stronger influence of omics-integrated features on the model output. Features such as `total_metabo_risk` and `total_transcript_risk` emerged as the most impactful drivers of metabolic risk predictions, reflecting the key role of molecular dysregulation in metabolic health stratification.

### D. Pathway Enrichment Analysis

To evaluate the biological relevance of omics-derived risk scores, we conducted pathway enrichment analysis across each of the five omics layers using exposure-weighted rankings. The findings are presented separately for prenatal and childhood stages, as the biological context and developmental timing strongly influence interpretation.

*1) Transcriptome:* **Prenatal Results:** GSEA on the ranked transcript-level data identified significant enrichment for pathways involved in B cell activation, G1/S cell cycle checkpoint control, and cytoplasmic translation. These findings suggest a strong engagement of immune system ontogeny and cell proliferation regulation during fetal development—mechanisms that may be sensitive to prenatal immunotoxic exposures.

**Childhood Results:** Enriched pathways included glutathione metabolism, T-cell receptor signaling, and programmed cell death. These reflect the sustained activation of antioxidant defenses and immune surveillance, indicating metabolic inflammation and homeostatic processes active in early childhood.
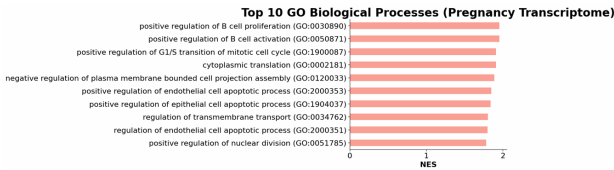

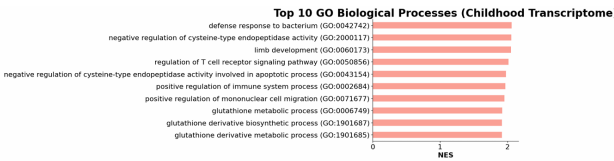Fig. 4. Transcriptome Pregnancy – Top GO Biological Processes


Fig. 5. Transcriptome Childhood – Top GO Biological Processes

*2) Methylome:* **Prenatal Results:** Epigenetically regulated genes were enriched for pathways involved in T-cell differentiation, epithelial cell polarity, and neurodevelopmental patterning. These findings point to long-term immune and tissue structuring effects rooted in fetal life.

**Childhood Results:** Key pathways included synaptic signaling, renal and reproductive system development, and embryonic morphogenesis, reflecting sustained epigenetic regulation of organ-specific functions beyond birth.
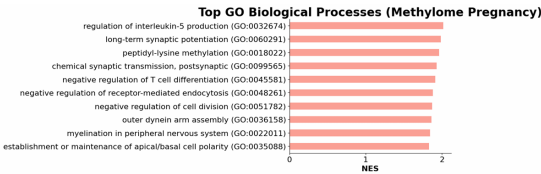

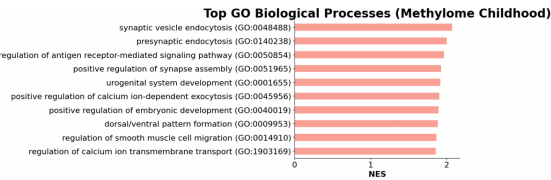Fig. 6. Methylome Pregnancy – Top GO Biological Processes


Fig. 7. Methylome Childhood – Top GO Biological Processes

*3) miRNA:* **Prenatal Results:** miEAA-based enrichment revealed dominant signals in vascular and tumor-related ontologies, including lung carcinoma, ischemic vascular disease, and liposarcoma. These suggest that even at the miRNA level, early exposure may influence vascular integrity and tumor suppressor regulation pathways.

**Childhood Results:** Enriched ontologies included immune and CNS conditions such as acute myeloid leukemia, cardiovascular disease, and multiple sclerosis. These may reflect systemic dysregulation of immune function and neuroinflammation—both relevant to pediatric metabolic risk.
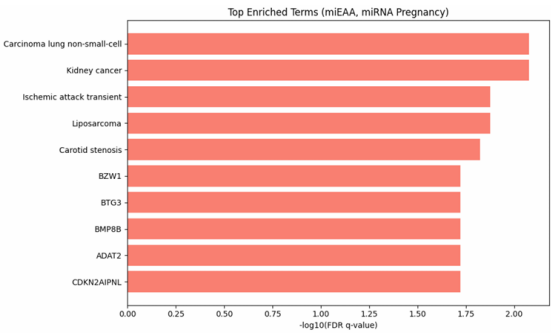

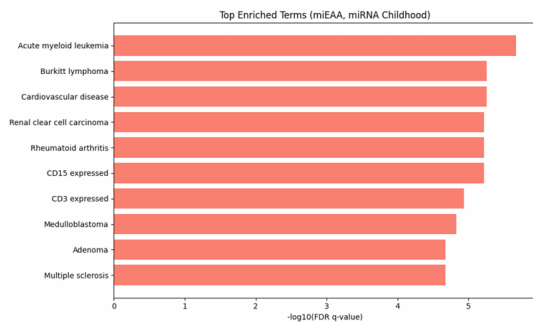Fig. 8. miRNA Pregnancy – Enriched Disease Ontologies

Fig. 9. miRNA Childhood – Enriched Disease Ontologies

*4) Proteome:* **Prenatal Results:** Protein-level exposure associations revealed enrichment for cytokine production, phosphorylation cascades, and cell proliferation—suggesting active immune modulation and intracellular signaling responses to in utero environmental stressors.

**Childhood Results:** Postnatal proteomic profiles were enriched for transcriptional regulation, biosynthetic pathways, and cell migration. This reflects ongoing tissue remodeling and energy regulation during early childhood growth.
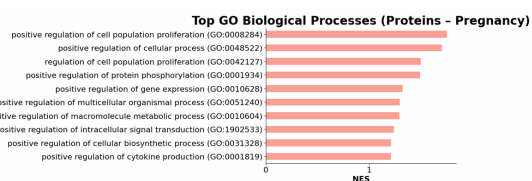


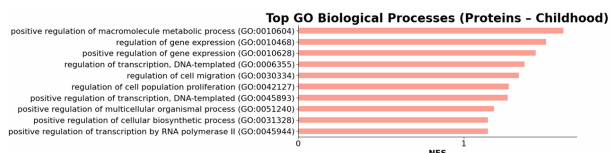Fig. 10. Proteome Pregnancy – Top GO Biological Processes



Fig. 11. Proteome Childhood – Top GO Biological Processes

*5) Metabolome:* **Results Across Both Periods:** Strikingly, both prenatal and childhood metabolomic profiles revealed nearly identical enrichment patterns, focused on amino acid metabolism, neurotransmitter transport (e.g., $Na^+/Cl^-$-dependent transporters), and membrane translocation systems. These pathways are central to metabolic homeostasis, cellular energy handling, and synaptic signaling.

*a) Integrated Interpretation and Clinical Relevance::* Across all five omics domains, the enrichment findings converge on four major biological themes: immune activation, oxidative stress regulation, biosynthetic remodeling, and neurodevelopmental signaling. These processes are well-established drivers of pediatric metabolic dysfunction and comorbidities such as obesity, insulin resistance, and cardiovascular disease.

What makes these findings clinically meaningful is their consistency across multiple molecular layers and develop-
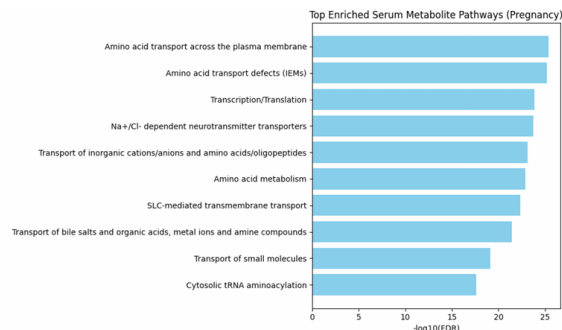


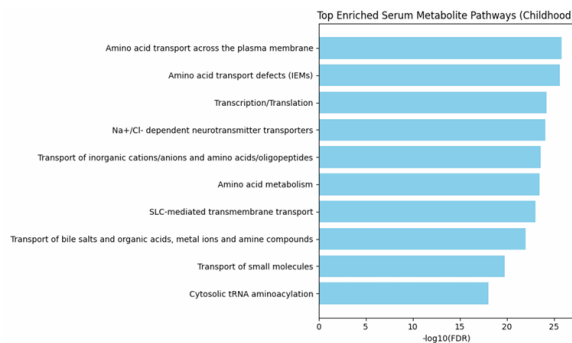Fig. 12. Metabolome Pregnancy – Enriched Pathways (SMPDB)



Fig. 13. Metabolome Childhood – Enriched Pathways (SMPDB)

mental windows. For example, the transcriptome and proteome both highlight immune cell signaling and antioxidant activity, while the methylome and miRNA profiles capture long-term regulation of neuroimmune systems. Meanwhile, the metabolome provides functional validation through altered energy substrates and amino acid transport pathways.

This cross-layer convergence confirms that the omics-derived risk scores do not merely reflect statistical noise or overfitting. Rather, they encode molecular signals that are functionally relevant, biologically interpretable, and temporally persistent. These signatures support the utility of omics-informed risk stratification not just for predictive modeling, but for identifying mechanistic intervention targets in pediatric metabolic health. In doing so, this work bridges predictive analytics with translational biology, reinforcing the scientific credibility of our framework for real-world clinical application.

*E. Overlap and Interpretability*

Several variables emerged as consistent predictors across both SHAP importance (from XGBoost) and LASSO feature selection, highlighting their robust influence on pediatric metabolic risk.

- `hs_pfna_c_Log2` – This represents logged PFNA levels measured in child serum. PFNA (Perfluorononanoic acid) is a long-chain perfluoroalkyl substance (PFAS) used in water- and stain-resistant materials. It is known to bioaccumulate and has been linked to lipid dysregulation, endocrine disruption, and early-life metabolic outcomes.

Its consistent selection indicates that PFNA exposure is a strong environmental determinant of metabolic risk in the HELIX cohort.

- `protein_child_risk` – This is the aggregated protein-based risk score calculated using exposure-weighted effect sizes. It reflects systemic changes in protein expression—particularly involving inflammatory cytokines, growth factors, or immune signaling proteins—induced by early-life environmental exposures. Its presence in both SHAP and LASSO suggests that protein-level immune and metabolic dysregulation plays a key role in shaping metabolic health trajectories.

- `h_PM_Log` – This variable captures long-term average PM2.5 or PM10 exposure, log-transformed for normalization. Airborne particulate matter is a well-documented cardiometabolic risk factor, capable of inducing oxidative stress, insulin resistance, and systemic inflammation. Its consistent inclusion reinforces the relevance of ambient air quality as a modifiable predictor of child health.

- `total_metabo_risk` – This score summarizes exposure-weighted metabolomic disruption. Metabolites reflect real-time biochemical states including energy metabolism, mitochondrial function, and oxidative stress. The prominence of this feature across both modeling approaches suggests that metabolic pathway alterations—especially those linked to lipid, amino acid, or glucose metabolism—are major mediators of early-life exposure effects on health.

The convergence of these features across orthogonal selection methods strengthens their interpretability and reinforces the biological plausibility of the model. Together, they span multiple biological layers—external exposure (PFNA, PM), proteomic signaling, and metabolic output—forming a coherent exposome-to-phenotype risk cascade.

## IV. DISCUSSION

This study presents a biologically grounded risk stratification framework that integrates environmental exposures with omics-informed molecular signatures. By combining exposure-weighted risk scores across five omics layers, we construct individualized metabolic risk profiles that capture both external insult and internal physiological response. Compared to exposure-only models, this integrated approach offers improved internal performance metrics and enhanced mechanistic interpretability. While these gains demonstrate the framework's capacity to consistently reproduce its stratification logic for unseen individuals within the HELIX dataset, they should not be interpreted as evidence of external predictive validity. Establishing such validation would require independent datasets containing both comparable omics measurements and clinically confirmed outcomes.

**Model Interpretation:** Convergent evidence from both SHAP and LASSO analyses highlights consistent high-impact features: PFNA exposure (`hs_pfna_c_Log2`), ambient particulate matter (`h_PM_Log`), and proteomic alterations (`protein_child_risk`). These features span the exposome-to-biology axis—linking persistent pollutants and air toxics to immune and metabolic dysregulation. The recurrence of `total_metabo_risk` further underscores the role of metabolomic disruption in mediating the downstream effects of early-life environmental exposures on pediatric health. These results demonstrate that biologically informed models not only predict well, but also reveal interpretable pathways of metabolic dysfunction with actionable insight.

**Biological Plausibility via Pathway Enrichment:** To confirm that omics-derived risk scores reflect meaningful biology rather than modeling artifacts, we performed pathway enrichment analysis across all five omics layers. Enriched pathways consistently mapped to domains implicated in pediatric metabolic disease, including immune activation (e.g., T-cell signaling, cytokine production), oxidative stress (e.g., glutathione metabolism), biosynthetic remodeling (e.g., protein phosphorylation, amino acid metabolism), and neurodevelopment (e.g., synaptic signaling, morphogenesis). These findings were consistent across both prenatal and childhood data, reinforcing the persistence of environmentally driven molecular disruption. The convergence of enriched biological processes with predictive features identified via SHAP and LASSO strengthens the clinical and mechanistic credibility of the stratification framework.

**Limitations:** The framework presented here does not aim to predict clinically diagnosed metabolic outcomes, but rather to construct a biologically informed composite score that reflects molecular susceptibility based on exposure–omics associations. Because the risk tiers are derived from the same variables used to compute the score, the reported performance measures primarily capture internal consistency rather than validation against an independent clinical endpoint. This inherent linkage limits direct assessment of predictive generalizability; however, the derived profiles closely align with well-established biological mechanisms—such as immune activation, oxidative stress, and endocrine disruption—known to contribute to early metabolic vulnerability. As such, the framework offers a valuable exploratory tool for identifying biologically plausible early warning signals that could inform prevention and intervention strategies unitl broader, outcome-linked datasets become available.

**Scientific Contribution:** We introduce a rule-based, exposure-weighted multi-omics modeling framework for stratifying pediatric metabolic risk that addresses a critical gap in the field: the absence of datasets containing both comprehensive omics measurements and validated clinical metabolic outcomes in children. Our pipeline is modular, interpretable, and biologically grounded—allowing integration of molecular response signatures with environmental burden to create mechanistically informed risk profiles. Given the current data landscape, we establish internal validation through pathway enrichment analysis, demonstrating that molecular features used in risk scoring correspond to well-established biological processes including immune activation, oxidative stress, and neurodevelopmental regulation—pathways directly implicated in metabolic dysfunction. This dual-layered inter-

pretability—SHAP for model behavior and pathway enrichment for functional coherence—provides robust evidence of biological plausibility while maintaining transparency about methodological constraints. As pediatric multi-omics datasets with clinical endpoints remain scarce, this work establishes a methodologically sound framework ready for deployment and validation when such ground truth data become available, thereby advancing the field's capacity for exposure-aware risk stratification in early childhood.

## V. Conclusion

Exposure-weighted omics risk profiles provide a biologically grounded approach to characterizing metabolic susceptibility in early childhood. By integrating systemic molecular alterations with detailed exposure histories, the combined framework demonstrates stronger performance than traditional exposure-only approaches. SHAP and LASSO analyses reveal convergent predictors—such as PFAS levels, particulate air pollution, and proteomic/metabolomic disruption—highlighting mechanistic pathways with translational relevance. Pathway enrichment further shows that top-ranked features correspond to immune, oxidative, neurodevelopmental, and biosynthetic processes implicated in early metabolic dysfunction. This work establishes a scalable and biologically interpretable foundation for exposure-aware risk stratification in pediatric health, offering clear potential for refinement and broader application as multi-omics resources continue to expand.

## Acknowledgments

## References

[1] Vaghari et al., "Multi-Omics Integration Approaches for Stratifying Disease Risk," Trends Mol Med., 2021.

[2] Lundberg and Lee, "A Unified Approach to Interpreting Model Predictions," NIPS, 2017.

[3] Tibshirani, "Regression Shrinkage and Selection via the LASSO," J. Royal Stat. Soc., 1996.

[4] Chen and Guestrin, "XGBoost: A Scalable Tree Boosting System," KDD, 2016.

[5] Maitre et al., "Environmental Influences on Child Health: The HELIX Study," Int. J. Epidemiol., 2018.

[6] Aguilar et al., "Integrative machine learning approaches for predicting disease risk using multi-omics data from the UK Biobank," *bioRxiv Preprint*, Apr.2024 – doi:10.1101/2024.04.16.589819 :contentReferenceindex=1

[7] Baião et al., "A technical review of multi-omics data integration methods: from classical statistical to deep generative approaches," *arXiv preprint*, Jan.292025 :contentReferenceindex=2

[8] Wang et al., "Emerging trends in systems biology: multi-omics integration and beyond," *Computational Molecular Biology*, vol.14, no.5, pp.211–219, 2024 – doi:10.5376/cmb.2024.14.0024 :contentReferenceindex=3

[9] MDPI et al., "Integrative multi-omics approach for stratification of tumor recurrence risk," *bioRxiv Preprint*, 2021 :contentReferenceindex=4

[10] Aget al., "Deep Learning-Based Multi-Omics Data Integration Reveals Two Neuroblastoma Subtypes," *Frontiers in Genetics*, 2018 :contentReferenceindex=5

[11] Translational Medicine et al., "Machine learning and multi-omics integration: advancing cardiovascular research," *J. Translational Medicine*, 2025; open access review :contentReferenceindex=6

[12] eLife, "Network-based multi-omics integration reveals at-risk metabolic profile in treated HIV patients," *eLife*, 2023 :contentReferenceindex=7

[13] MDPI, "Navigating Challenges and Opportunities in Multi-Omics Integration: Wearable & circadian data," *Medicines*, 2024 :contentReferenceindex=8

[14] Zhang et al., "OmiEmbed: a unified multi-task deep learning framework for multi-omics data," *arXiv preprint*, Feb.2021 :contentReferenceindex=9

[15] Ma et al., "Integrate Any Omics: Towards genome-wide data integration for patient stratification," *arXiv preprint*, Jan.2024 :contentReferenceindex=10

[16] Alharbi et al., "Comparative analysis of multi-omics integration using graph neural networks for cancer classification," *arXiv preprint*, Oct.2024 :contentReferenceindex=11

[17] Mankoo et al., "Integrated prediction of ovarian cancer recurrence using gene, miRNA, methylation, and CNA," *PMC Open Access Article*, year2022/2023 :contentReferenceindex=12

[18] "Applications of multi-omics analysis in human diseases: filtering novel biomolecule-phenotype associations," *PMC Review*, 2023 :contentReferenceindex=13

[19] Shetty et al., "Network-based integrative multi-omics approach reveals COVID-19 disease-state specific biosignatures and interactions," *Frontiers in Molecular Biosciences*, 2024 :contentReferenceindex=14

[20] "Multi-omics approaches to disease," Genome Biology, 2017 – review on combining omics to stratify risk :contentReferenceindex=15

[21] Bersanelli et al., "Methods for the integration of multi-omics data: mathematical aspects," *BMC Bioinformatics*, 2016 :contentReferenceindex=16

[22] Argelaguet et al., "MOFA+: a statistical framework for comprehensive integration of multi-modal single-cell data," *Genome Biology*, 2020 :contentReferenceindex=17

[23] Stuart Butler et al., "Integrated analysis of multimodal single-cell data," *Cell*, 2021 :contentReferenceindex=18

[24] Yuan et al., "Advances in bulk and single-cell multi-omics approaches for systems biology," *Briefings in Bioinformatics*, 2021 :contentReferenceindex=19

[25] Magnuson et al., "Latent embeddings for predictive modeling in multi-omics," *Nat Rev Nephrol*, 2021 review :contentReferenceindex=20

[26] Conesa et al., "PaintOmics 3: web resource for pathway analysis and visualization of multi-omics data," *Nucleic Acids Research*, 2018 :contentReferenceindex=21

[27] Suo et al., "Deep learning-based integration predicts survival in liver cancer," *Cancer Research*, 2018 :contentReferenceindex=22

[28] Chen et al., "A multivariate approach to integration of multi-omics datasets," *BMC Bioinformatics*, 2014 :contentReferenceindex=23

[29] Nicora et al., "Network approaches to systems biology analysis of complex disease," *Briefings in Bioinformatics*, 2017 :contentReferenceindex=24

[30] Hasin et al., "Multi-omics approaches to disease," *Genome Biology*, 2017 review :contentReferenceindex=25