

Lecture 12:

# **Multi-Modal Data Integration**

Integrative biology

Systems medicine

Holistic approaches

**Introduction to Biomedical Data Science**

# Lecture Contents

**Part 1:** Integration Methods

**Part 2:** Multi-Omics Applications

**Part 3:** Clinical Applications and Future Directions

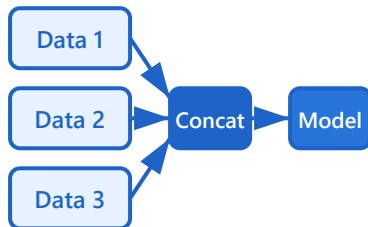
**Part 1/3:**

# Integration Methods

- Mathematical frameworks
- Computational approaches
- Evaluation metrics

## Early vs Late Fusion

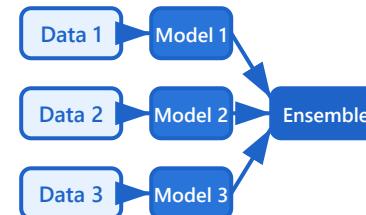
### Early Fusion



Feature-level concatenation

- ✓ Captures feature interactions
- ✓ Joint representation learning
- ✗ Computationally expensive

### Late Fusion



Decision-level combination

- ✓ Flexible and modular
- ✓ Independent training
- ✗ Misses feature interactions

### Intermediate Fusion

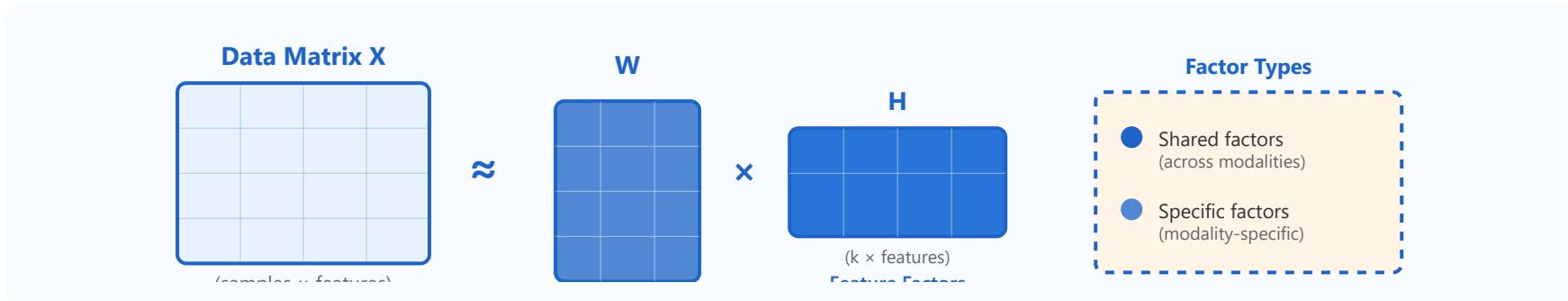


Combines features at intermediate layers, balancing both approaches

### Trade-offs

Early: captures interactions but computationally expensive | Late: flexible but misses feature interactions

# Matrix Factorization for Multi-Omics Integration



## NMF Methods

Non-negative Matrix Factorization for parts-based representation

## Joint NMF

Simultaneous factorization of multiple data matrices

## iCluster

Integrative clustering of multiple cancer genomic data types

## Integrative NMF

Shared and data-specific factors

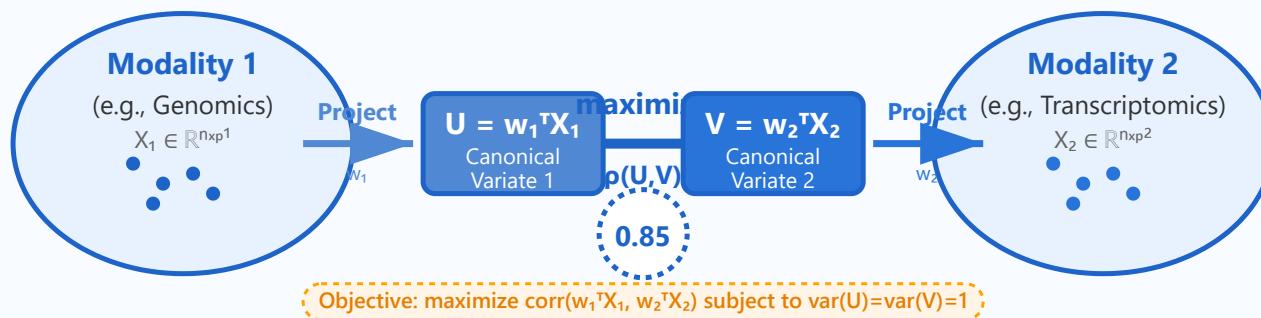
## Interpretation

Biological meaning of latent factors

# Canonical Correlation Analysis (CCA)

**Canonical Correlation Analysis (CCA)** is a multivariate statistical method that explores the relationships between two sets of variables. It finds linear combinations of variables in each set that have maximum correlation with each other. This powerful technique is widely used in genomics, neuroscience, psychology, and machine learning for integrating multi-modal data.

The fundamental goal is to identify patterns of association between two data matrices  $X_1$  ( $n \times p_1$ ) and  $X_2$  ( $n \times p_2$ ), where  $n$  is the number of samples and  $p_1, p_2$  are the dimensions of each modality.



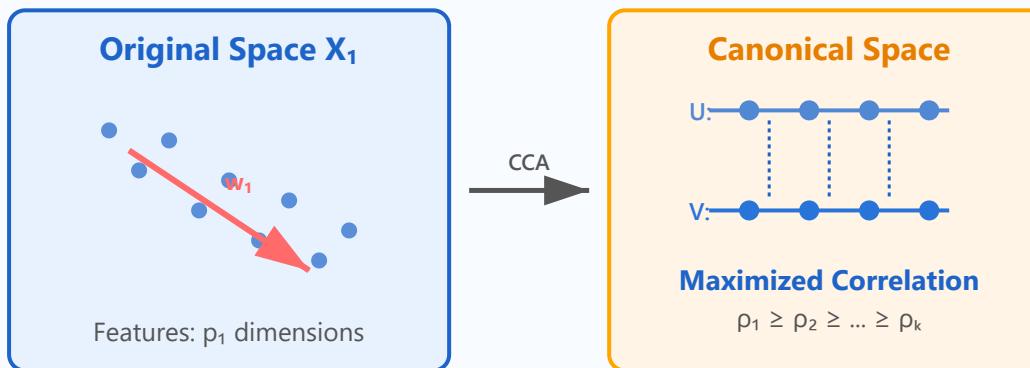
## Detailed Methods and Applications

## 1

# Classical CCA Principles

Classical Canonical Correlation Analysis seeks to find linear combinations of variables from two datasets that are maximally correlated. The method identifies canonical weight vectors  $w_1$  and  $w_2$  that maximize the correlation between the projected variables.

```
maximize ρ = corr(w1TX1, w2TX2)
subject to: var(w1TX1) = var(w2TX2) = 1
```



## Key Characteristics:

- **Multiple canonical correlations:** CCA finds  $k$  pairs of canonical variates ( $k = \min(p_1, p_2)$ ), ordered by decreasing correlation
- **Orthogonality:** Successive canonical variates are uncorrelated with previous ones
- **Dimensionality reduction:** Projects high-dimensional data into lower-dimensional canonical space
- **Mathematical solution:** Solved via generalized eigenvalue decomposition

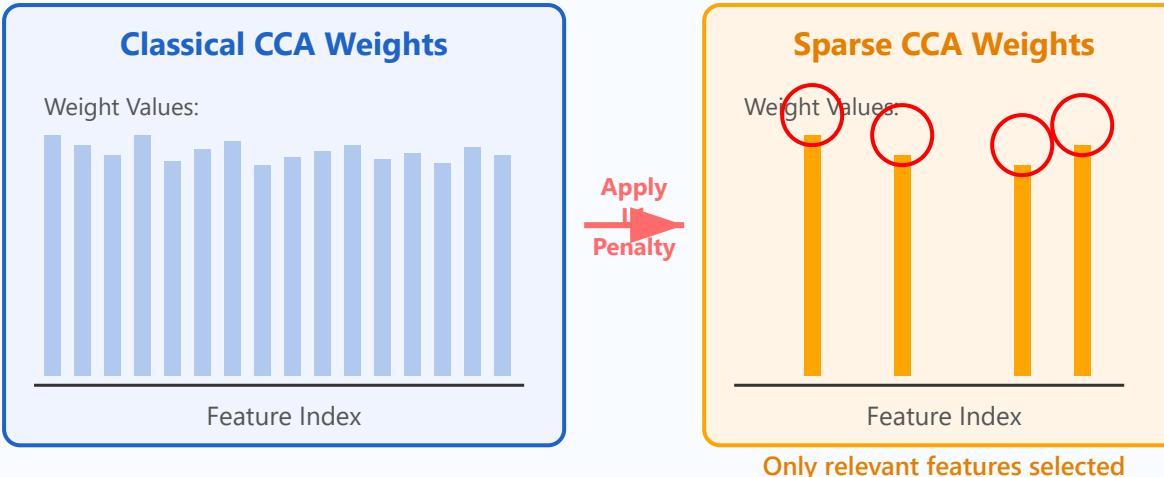
### Real-world Applications:

- **Genomics:** Correlating gene expression with clinical outcomes
- **Neuroscience:** Linking brain imaging with behavioral measures
- **Economics:** Relating economic indicators across countries
- **Psychology:** Connecting psychological test scores with behavioral assessments

## 2 Sparse CCA

Sparse Canonical Correlation Analysis addresses the challenge of high-dimensional data by introducing sparsity constraints on the canonical weight vectors. This approach performs simultaneous feature selection and correlation maximization, making results more interpretable and identifying the most relevant features.

```
maximize ρ = corr(w1TX1, w2TX2)
subject to: ||w1||1 ≤ c1, ||w2||1 ≤ c2 (L1 penalty for sparsity)
```



### Key Characteristics:

- **Feature selection:** Automatically identifies the most important variables in each modality
- **Interpretability:** Sparse weight vectors are easier to interpret than dense ones
- **Regularization methods:** L1 (LASSO), elastic net, or structured sparsity penalties
- **Computational efficiency:** Reduces computational burden in high-dimensional settings
- **Overfitting prevention:** Reduces model complexity and improves generalization

### Real-world Applications:

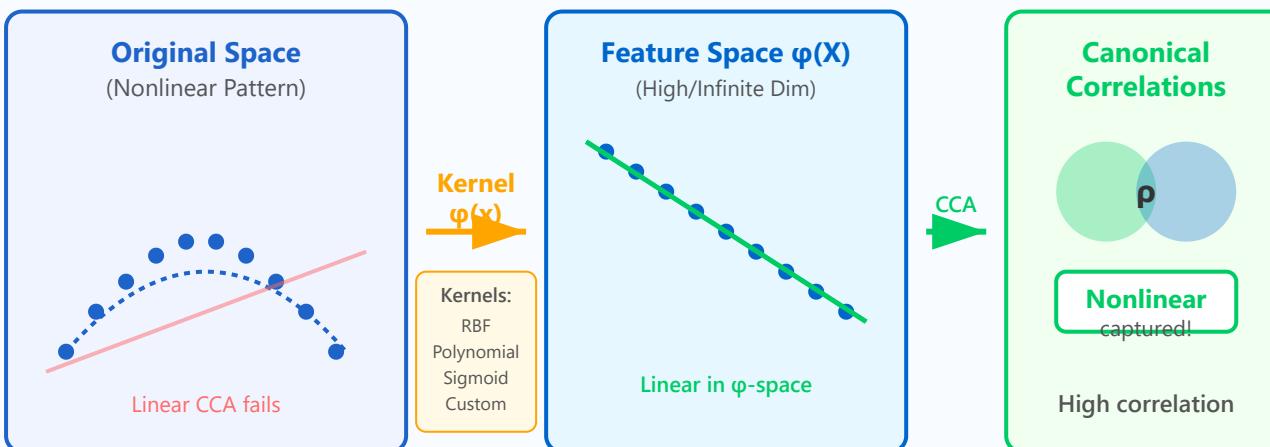
- **Genomics:** Identifying key genes associated with phenotypes from thousands of candidates
- **Medical imaging:** Selecting relevant brain regions correlated with cognitive scores

- **Bioinformatics:** Finding biomarkers linking proteomics and metabolomics data
- **Text analysis:** Identifying important words linking document collections

### 3 Kernel CCA

Kernel Canonical Correlation Analysis extends classical CCA to capture nonlinear relationships between datasets. By mapping data into high-dimensional (potentially infinite) feature spaces using kernel functions, KCCA can identify complex, nonlinear associations that linear CCA would miss.

```
maximize ρ = corr(α1TK1, α2TK2)
where K1, K2 are kernel matrices: K(x, x') = φ(x)Tφ(x')
```



## Key Characteristics:

- **Nonlinear relationships:** Captures complex, nonlinear associations between data modalities
- **Kernel trick:** Avoids explicit computation in high-dimensional space using  $K(x,x') = \varphi(x)^T\varphi(x')$
- **Popular kernels:** RBF/Gaussian ( $\exp(-\gamma||x-x'||^2)$ ), polynomial  $((x\cdot x' + c)^d)$ , sigmoid
- **Regularization:** Often requires regularization to prevent overfitting in feature space
- **Computational cost:**  $O(n^3)$  complexity due to kernel matrix operations

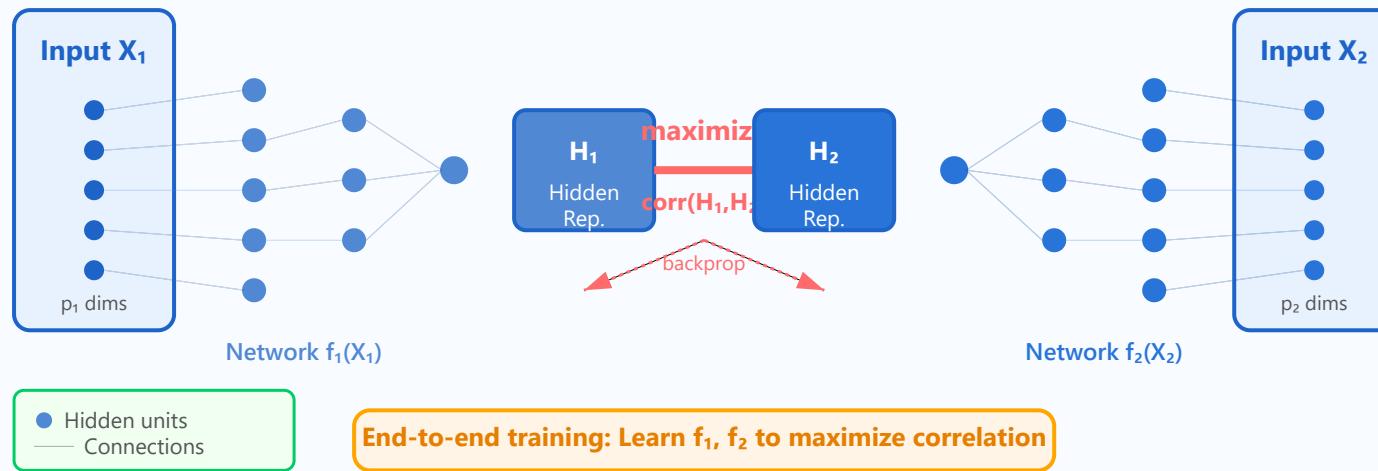
## Real-world Applications:

- **Computer vision:** Matching images and text descriptions with complex visual features
- **Bioinformatics:** Capturing nonlinear gene-phenotype relationships
- **Speech recognition:** Correlating acoustic features with linguistic representations
- **Climate science:** Linking atmospheric variables with nonlinear interactions

## 4 Deep CCA

Deep Canonical Correlation Analysis leverages deep neural networks to learn nonlinear transformations of the input data that maximize correlation. Unlike Kernel CCA which uses fixed transformations, Deep CCA learns optimal representations through end-to-end training, making it highly flexible and powerful for complex, high-dimensional data.

```
maximize corr(f1(X1; θ1), f2(X2; θ2))  
where f1, f2 are deep neural networks with parameters θ1, θ2
```



### Key Characteristics:

- **Learned transformations:** Neural networks learn optimal nonlinear mappings instead of using predefined kernels
- **Scalability:** Can handle very high-dimensional inputs more efficiently than Kernel CCA
- **End-to-end training:** Networks are trained jointly using gradient-based optimization (backpropagation)
- **Flexibility:** Network architectures can be customized (CNNs for images, RNNs for sequences, etc.)
- **Modern variants:** DCCA, DCCAE (with autoencoders), Variational Deep CCA

### Real-world Applications:

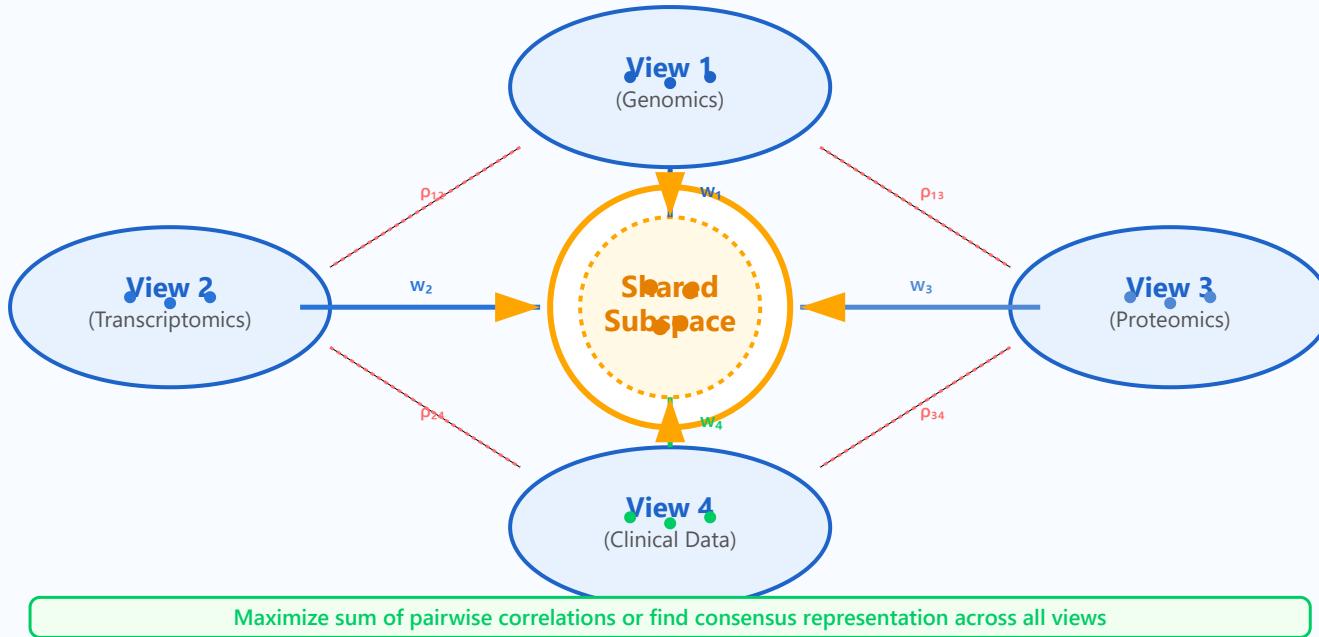
- **Multi-modal learning:** Correlating images with captions, videos with audio
- **Cross-lingual NLP:** Learning shared representations across languages

- **Medical imaging:** Integrating MRI, CT, and PET scans with clinical data
- **Sensor fusion:** Combining data from multiple sensors in robotics and IoT

## 5 Multi-view CCA (Generalized CCA)

Multi-view Canonical Correlation Analysis extends classical CCA beyond two datasets to handle multiple (more than two) views or modalities simultaneously. This is essential in modern data integration problems where information comes from multiple sources that need to be jointly analyzed.

```
maximize Σi,j w(i,j) · corr(wiTXi, wjTXj)
where i, j ∈ {1, 2, ..., M} and M > 2 (multiple views)
```



### Key Characteristics:

- **Multiple views:** Handles  $M > 2$  data modalities simultaneously (e.g., genomics, transcriptomics, proteomics, clinical data)
- **Consensus representation:** Finds a shared subspace that captures common information across all views
- **Different objectives:** Sum of pairwise correlations (SUMCOR), maximum variance (MAXVAR), or generalized eigenvalue problem
- **Incomplete views:** Can handle missing data in some views for some samples
- **Weighting schemes:** Can assign different importance weights to different view pairs

#### Popular variants:

- SUMCOR:  $\text{maximize } \sum_{i < j} \text{corr}(U_i, U_j)$
- MAXVAR:  $\text{maximize } \sum_i \text{var}(U_i) \text{ subject to consensus constraint}$
- Generalized CCA: eigen-decomposition of cross-covariance matrices

### Real-world Applications:

- **Multi-omics integration:** Combining genomics, transcriptomics, proteomics, and metabolomics data
- **Multi-modal medical imaging:** Integrating MRI, CT, PET, and ultrasound with clinical records
- **Social media analysis:** Analyzing text, images, and user behavior across platforms
- **Climate modeling:** Integrating satellite data, ground sensors, and simulation outputs
- **Recommender systems:** Combining user behavior, demographics, content features, and social networks

## Comparison Summary

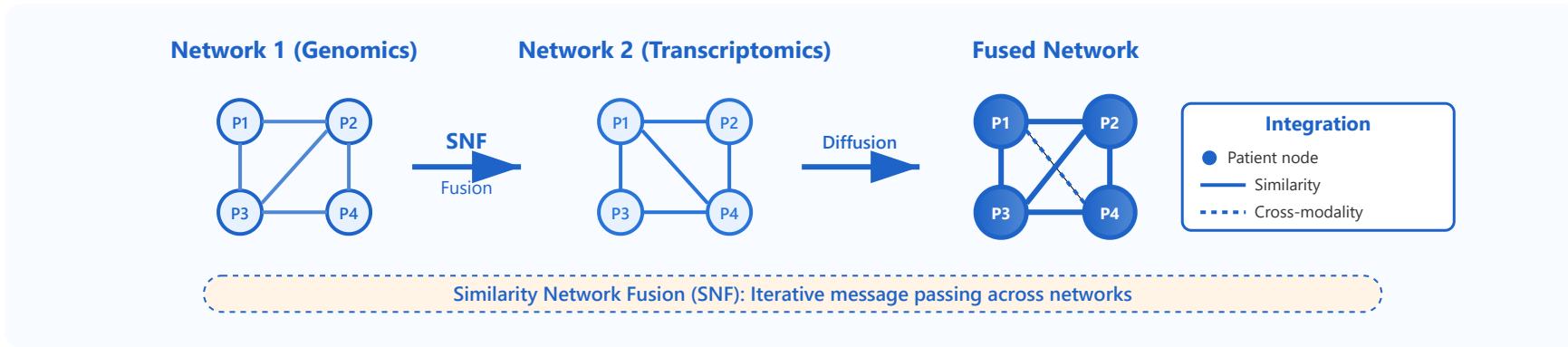
Method	Linearity	Sparsity	# Views	Scalability
Classical CCA	Linear	Dense	2	Medium ( $O(p^3)$ )
Sparse CCA	Linear	<span style="color: green;">● Sparse</span>	2	Good (feature selection)
Kernel CCA	<span style="color: red;">● Nonlinear</span>	Dense	2	Poor ( $O(n^3)$ )
Deep CCA	<span style="color: red;">● Nonlinear (learned)</span>	Flexible	2	<span style="color: green;">● Excellent</span>
Multi-view CCA	Linear (typically)	Varies	<span style="color: red;">● <math>M &gt; 2</math></span>	Depends on M

### Choosing the Right Method:

- **Classical CCA:** Best for moderate dimensions, linear relationships, interpretable results
- **Sparse CCA:** High-dimensional data where feature selection is crucial (e.g., genomics)

- **Kernel CCA:** Strong nonlinear relationships with moderate sample sizes
- **Deep CCA:** Complex nonlinear patterns, very high dimensions, large datasets available
- **Multi-view CCA:** More than two data modalities need to be integrated simultaneously

## Graph-based Integration



### Similarity Networks

Patient or feature similarity graphs

### Network Fusion

SNF: fusing multiple similarity networks

### Random Walk

Diffusion-based integration on networks

### Graph Neural Networks

Deep learning on graph-structured data

### Multiplex Networks

Multi-layer network representations

## 1. Similarity Networks

## Overview

Similarity networks represent relationships between patients or features based on their shared characteristics across omics datasets. Each node represents a sample (patient) or feature (gene, protein), and edges represent similarity scores.

## Construction Methods

- **K-Nearest Neighbors (KNN):** Connect each node to its k most similar neighbors
- **Gaussian Kernel:** Weight edges using exponential decay based on distance
- **Correlation-based:** Use Pearson or Spearman correlation coefficients

## Applications

- Patient stratification and subtype discovery
- Feature selection and biomarker identification
- Disease trajectory modeling

### Key Points:

- Edge weights represent similarity strength between patients
- Thresholding removes weak connections for clearer network structure
- Multiple similarity metrics can be combined for robust networks

## Patient Similarity Network Construction

Gene Expression Data

P1	P2	P3	P4

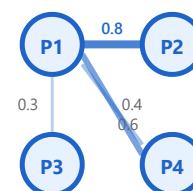
Calculate Similarity

Similarity Matrix

1.0	0.8	0.3	0.5
0.8	1.0	0.4	0.6

Build Network

Similarity Network



## 2. Network Fusion (SNF)

### Similarity Network Fusion Algorithm

SNF integrates multiple patient similarity networks by iteratively updating each network using information from other networks. This cross-network information exchange creates a unified representation that captures complementary information.

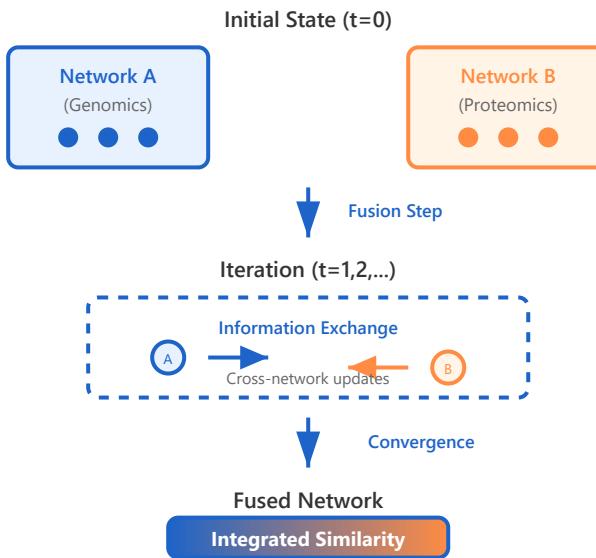
### Algorithm Steps

- **Step 1:** Construct individual similarity networks for each data type
- **Step 2:** Normalize networks to make them comparable
- **Step 3:** Iteratively update each network using local and global information
- **Step 4:** Converge to a fused network representation

### Advantages

- Handles incomplete data naturally
- Emphasizes concordant patterns across data types
- No feature-level alignment required

### SNF Iterative Process



### Mathematical Formulation:

- $P(t+1) = S \times [P(t) + \Sigma(P_k(t))] / (m-1) \times S^T$
- S: Local similarity matrix (KNN structure)
- P: Full similarity matrix being updated
- Converges typically within 10-20 iterations

### 3. Random Walk with Restart (RWR)

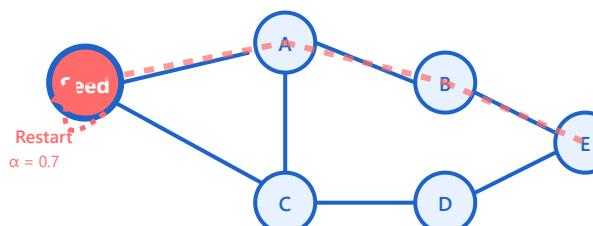
#### Concept

Random walk algorithms simulate a particle moving randomly across network edges, with a probability of restarting at seed nodes. This diffusion process captures both direct and indirect relationships in the network.

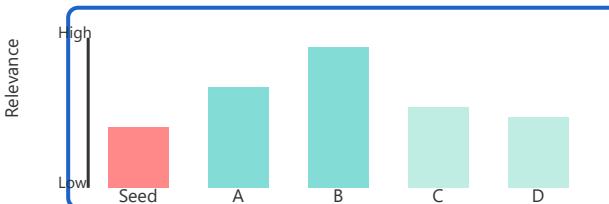
#### Algorithm

- **Initialization:** Start at seed node(s) representing genes, patients, or features of interest
- **Walk:** Move to neighboring nodes with probability proportional to edge weights
- **Restart:** With probability  $\alpha$  (typically 0.7), return to seed nodes
- **Convergence:** Iterate until steady-state probabilities stabilize

**Random Walk with Restart**



**Steady-State Probability**



## Applications in Multi-Omics

- Disease gene prioritization across omics layers
- Drug-target prediction in biological networks
- Pathway enrichment analysis

### Key Advantages:

- Captures global network topology, not just local neighborhoods
- Naturally handles heterogeneous networks with multiple node types
- Robust to noise and missing edges in biological networks
- Steady-state probabilities provide quantitative ranking of nodes

## 4. Graph Neural Networks (GNNs)

### Overview

Graph Neural Networks extend deep learning to graph-structured data, learning node representations by aggregating information from neighboring nodes. GNNs are particularly powerful for multi-omics integration as they can model complex relationships.

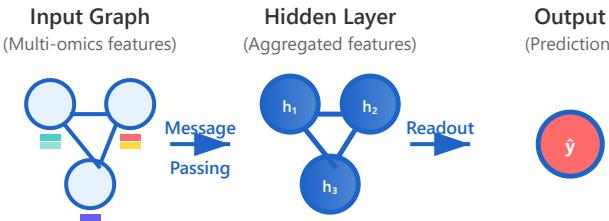
### Key Architectures

- **Graph Convolutional Networks (GCN):**  
Aggregate neighbor features using spectral convolutions
- **Graph Attention Networks (GAT):** Learn importance weights for each neighbor
- **GraphSAGE:** Sample and aggregate from fixed-size neighborhoods
- **Message Passing Neural Networks:** Generalized framework for node updates

## Multi-Omics Applications

- Patient outcome prediction using heterogeneous networks
- Drug response modeling with gene-gene interactions
- Biological pathway discovery

## Graph Neural Network Architecture



### Message Passing Mechanism

1. Aggregate  
Collect neighbor features
2. Combine  
Mix with own features
3. Update  
Apply neural network

#### Mathematical Form:

$$h'_{i,j} = \sigma(W \cdot \sum_j (a_{ij} \cdot h_j) + b)$$

$a_{ij}$ : attention weight,  $\sigma$ : activation function

### GNN Advantages for Multi-Omics:

- End-to-end learning of representations from raw data
- Automatic feature extraction from graph structure
- Can handle heterogeneous node and edge types
- Scalable to large biological networks

## 5. Multiplex Networks

### Concept

Multiplex networks represent multi-omics data as multiple layers where each layer corresponds to a different data type or interaction type. Nodes are shared across layers, but edges differ, capturing complementary relationships.

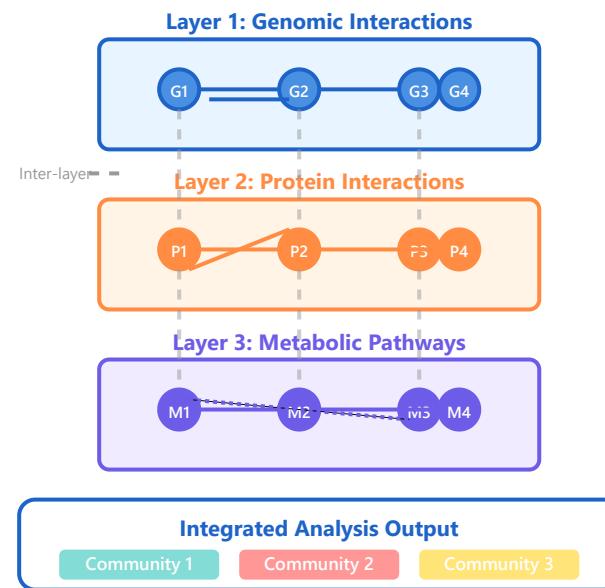
### Layer Types in Multi-Omics

- **Genomic Layer:** Gene-gene interactions, co-expression
- **Protein Layer:** Protein-protein interactions
- **Metabolic Layer:** Metabolite-metabolite relationships
- **Regulatory Layer:** TF-gene, miRNA-gene interactions

### Analysis Methods

- Multi-layer centrality measures
- Cross-layer community detection
- Multiplex random walks
- Layer coupling analysis

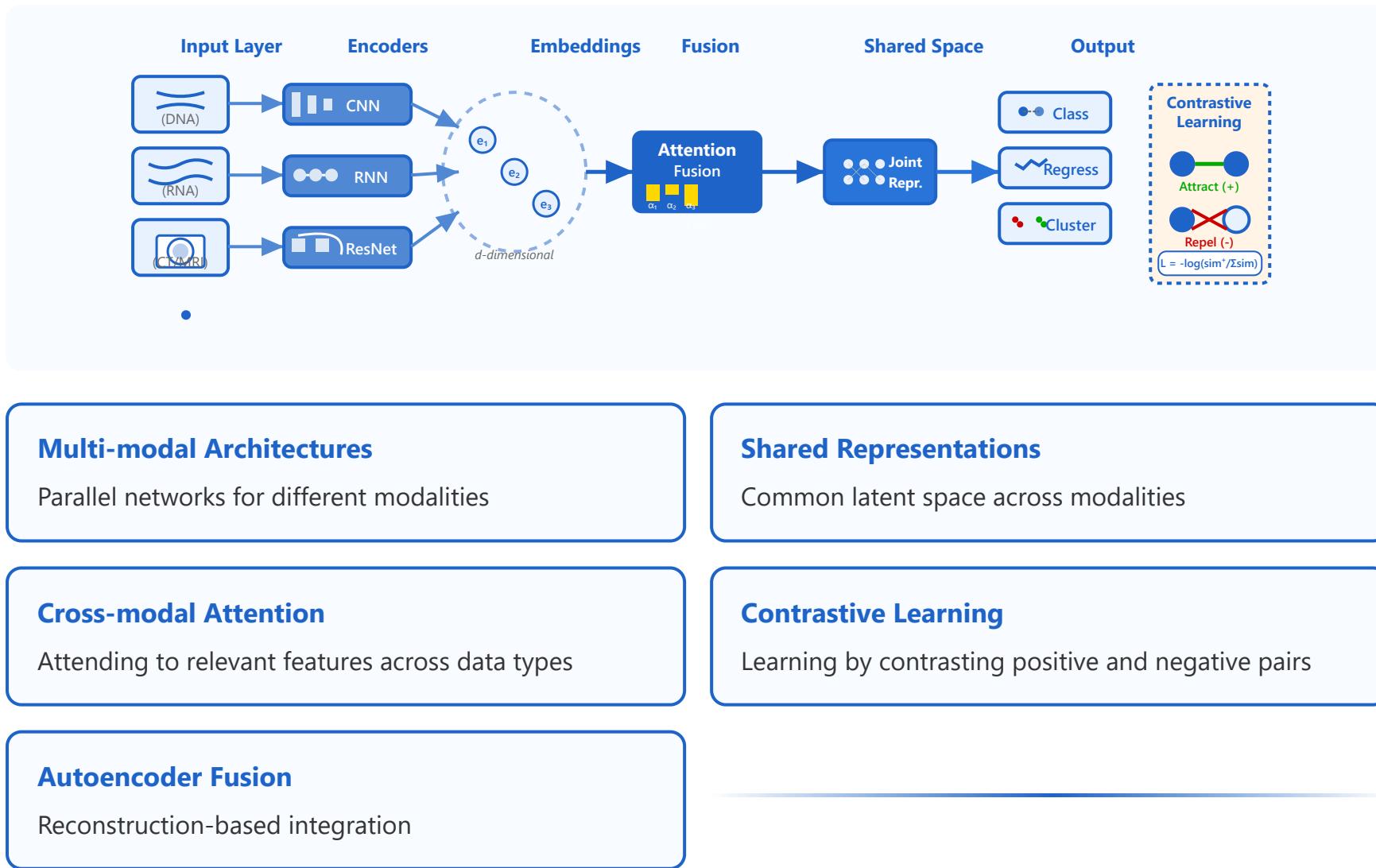
### Multiplex Network Structure



### **Key Features of Multiplex Networks:**

- Each layer captures a different biological relationship type
- Nodes represent the same biological entities across layers
- Inter-layer connections represent cross-omics relationships
- Enables discovery of layer-specific and cross-layer patterns
- Useful for identifying versatile biomarkers active across multiple omics layers

# Deep Learning Fusion Strategies



## Detailed Explanations and Examples

### 1 Multi-modal Architectures

In-depth analysis of each fusion strategy with visual examples

Multi-modal architectures employ parallel neural networks, each specialized for processing a specific data modality. This approach recognizes that different data types such as images, text, genomic sequences, and clinical measurements have unique structural characteristics that require specialized feature extraction methods.

## Architecture Design

The core principle is to use modality-specific encoders that transform raw input data into meaningful latent representations. For example, Convolutional Neural Networks (CNNs) excel at processing spatial data like medical images, Recurrent Neural Networks (RNNs) or Transformers handle sequential data like genomic sequences, and fully connected networks process tabular clinical data.

### Key Characteristics:

- **Modality-Specific Processing:** Each encoder is designed to capture the unique patterns and structures inherent to its input modality
- **Independent Feature Extraction:** Encoders operate independently before fusion, allowing parallel processing and specialized optimization

- **Flexible Integration:** Different fusion strategies can be applied (concatenation, addition, attention) to combine learned representations
- **Scalability:** New modalities can be added by simply incorporating additional encoder branches

#### Real-World Example: Cancer Diagnosis

A cancer diagnosis system might use a ResNet-50 CNN to process histopathology images, a Transformer to analyze gene expression profiles, and a fully connected network to process patient clinical data including age, biomarkers, and medical history. Each encoder extracts complementary information that is then fused for the final diagnosis.

## 2 Shared Representations

Shared representation learning aims to project data from different modalities into a common latent space where semantic relationships are preserved. This enables the model to learn unified

## 3 Cross-modal Attention

Cross-modal attention mechanisms allow the model to selectively focus on relevant features from one modality when processing another. This dynamic weighting mechanism learns which parts of each

representations that capture the underlying structure shared across modalities, facilitating cross-modal understanding and retrieval.

## Mathematical Formulation

$$z = f(x_1, x_2, \dots, x_n) \text{ where } z \in \mathbb{R}^d$$

(shared space)

The goal is to learn projection functions that map different modalities into the same dimensional space while preserving semantic similarity. Data points that are semantically similar should be close in this shared space, regardless of their original modality.

### Key Characteristics:

- **Semantic Alignment:** Points with similar meanings cluster together regardless of modality
- **Cross-modal Retrieval:** Enables finding relevant content across modalities through nearest neighbor search
- **Dimensionality Reduction:** Projects high-dimensional heterogeneous data into a unified lower-dimensional space
- **Transfer Learning:** Knowledge learned from one modality can benefit others through the shared representation

modality are most informative for the task at hand, enabling sophisticated feature interaction and complementary information extraction.

## Attention Mechanism

$$\text{Attention}(Q, K, V) = \text{softmax}(QK^T/\sqrt{d_k})V$$

In cross-modal attention, queries come from one modality while keys and values come from another. The attention weights determine how much each element of the source modality should contribute to the representation of the target modality.

### Key Characteristics:

- **Selective Focus:** Dynamically determines which features from source modality are relevant for target modality
- **Learnable Weights:** Attention weights are learned during training based on feature compatibility
- **Context-Dependent:** Attention patterns change based on the input, allowing adaptive information flow
- **Bidirectional:** Can be applied in both directions allowing mutual enhancement between modalities

## Types of Cross-modal Attention

## Training Objectives

Common training objectives include minimizing the distance between corresponding samples from different modalities, canonical correlation analysis (CCA) to maximize correlation, and triplet loss to ensure similar samples are closer than dissimilar ones.

### Real-World Example: Medical Report Generation

In radiology, X-ray images and diagnostic reports can be projected into a shared space. Given a new X-ray image, the system can retrieve similar images or generate relevant text descriptions by finding nearby points in the shared space. This enables automated report generation and similar case retrieval.

**Co-attention:** Both modalities attend to each other simultaneously, creating bidirectional information flow. **Self-attention across modalities:** Treats concatenated multi-modal features as a sequence and applies self-attention. **Guided attention:** One modality acts as a guide to selectively extract information from another.

### Real-World Example: Visual Question Answering

When answering "What color is the car?", the text encoder processes the question to generate queries, while the image encoder provides keys and values from different image regions. The attention mechanism focuses on regions containing cars, particularly where color information is visible, effectively grounding the linguistic query in visual content.

4

## Contrastive Learning

Contrastive learning trains models by learning to distinguish between similar (positive) and dissimilar

5

## Autoencoder Fusion

Autoencoder-based fusion approaches learn multi-modal representations by training the model to

(negative) pairs of data. In multi-modal contexts, positive pairs consist of corresponding data from different modalities, such as an image and its caption, while negative pairs are non-matching combinations. The model learns representations where positive pairs are pulled together in the embedding space while negative pairs are pushed apart.

## Loss Function

$$L = -\log \left( \frac{\exp(\text{sim}(z_i, z_j)/\tau)}{\sum_k \exp(\text{sim}(z_i, z_k)/\tau)} \right)$$

Where  $\text{sim}(\cdot, \cdot)$  is a similarity function like cosine similarity,  $\tau$  is a temperature parameter, and the sum is over all negative samples. This is known as the InfoNCE loss.

### Key Characteristics:

- **Self-Supervised:** Doesn't require manual labels, uses the natural correspondence between modalities
- **Scalability:** Can leverage large amounts of unlabeled multi-modal data from the internet
- **Discriminative Learning:** Learns by comparison rather than reconstruction, focusing on distinctive features

reconstruct inputs from a fused latent representation. The key insight is that to successfully reconstruct multiple modalities from a shared bottleneck representation, the model must learn to capture the essential complementary information from all sources. This reconstruction objective naturally encourages the learning of comprehensive multi-modal features.

## Architecture Components

The architecture consists of modality-specific encoders that compress each input into latent representations, a fusion module that combines these representations (via concatenation, addition, or more complex operations), and modality-specific decoders that attempt to reconstruct the original inputs from the fused representation.

### Key Characteristics:

- **Unsupervised Learning:** Uses reconstruction as a self-supervised signal without requiring labeled data
- **Information Bottleneck:** The fusion layer acts as a bottleneck forcing compression of essential multi-modal information
- **Completeness:** Successful reconstruction requires capturing complementary information from all modalities

- **Transfer Learning:** Representations learned via contrastive learning transfer well to downstream tasks

- **Regularization:** Can add additional constraints like sparsity or disentanglement to the latent space

## Popular Frameworks

### **CLIP (Contrastive Language-Image Pre-training):**

Trains image and text encoders jointly by maximizing similarity between matching image-text pairs. **SimCLR:** Creates positive pairs through data augmentation of the same sample. **MoCo**

**(Momentum Contrast):** Uses a momentum encoder and queue to maintain a large number of negative samples.

### Real-World Example: Zero-Shot Classification

CLIP, trained on 400 million image-text pairs, can classify images into categories it has never explicitly seen during training. Given an image and a set of text descriptions like "a photo of a cat", "a photo of a dog", the model computes similarities between the image embedding and each text embedding, choosing the highest scoring match. This works because contrastive learning created a shared semantic space.

## Variants and Extensions

**Variational Autoencoders (VAE):** Add probabilistic modeling to the latent space, learning distributions rather than point estimates. This enables generation of new samples and better uncertainty quantification. **Multi-modal VAE:** Extensions that can handle missing modalities during inference by marginalizing over the latent distributions. **Cross-modal Autoencoders:** Train to reconstruct one modality from another, learning cross-modal mappings.

## Training Considerations

Reconstruction losses for different modalities may have different scales, requiring careful weighting. Common approaches include normalizing losses, using adaptive weighting schemes, or employing uncertainty-based weighting where the model learns optimal loss weights during training. Additionally, pre-training encoders separately before fusion can improve convergence.

### Real-World Example: Multi-omics Data Integration

In cancer research, scientists integrate genomics, transcriptomics, and proteomics data using autoencoder fusion. Each omics layer is encoded into a latent representation, fused in a bottleneck layer, then reconstructed. The fused representation captures the essential biological state, enabling patient stratification and biomarker discovery. The reconstruction objective ensures that no critical information from any single omics layer is lost, while the bottleneck forces the model to learn the most informative integrated features.

## Summary and Comparison

Each fusion strategy offers unique advantages suited to different scenarios. The choice depends on factors including data characteristics, computational resources, availability of labels, and the specific downstream task.

### Selection Guidelines

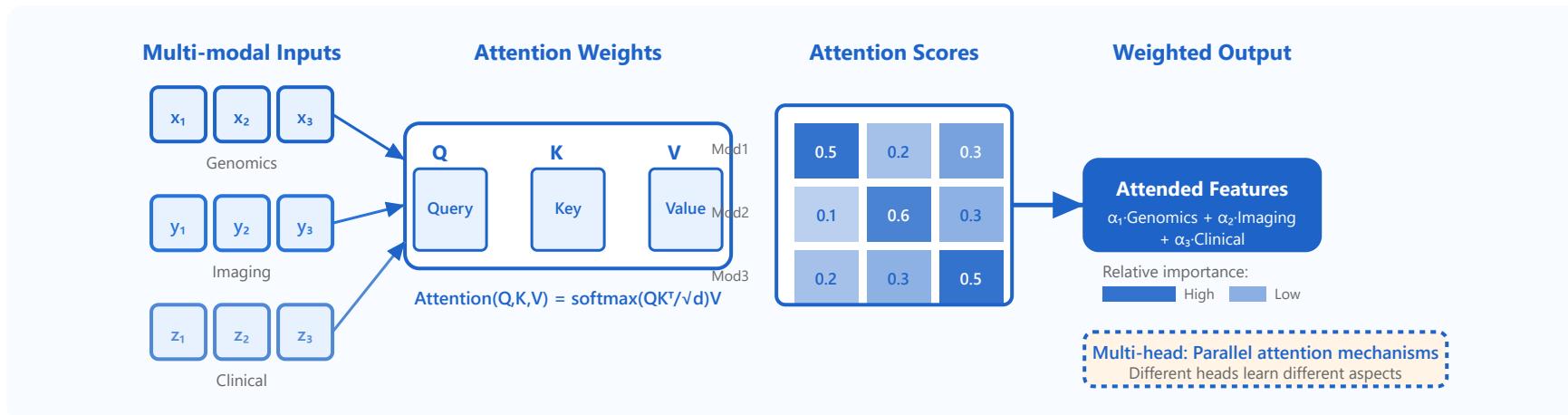
- **For Supervised Tasks with Labels:** Multi-modal architectures or cross-modal attention provide direct optimization for the target task
- **For Large Unlabeled Datasets:** Contrastive learning excels at leveraging web-scale data without manual annotation
- **For Cross-modal Retrieval:** Shared representations create a unified space enabling efficient similarity search
- **For Incomplete Data:** Autoencoder fusion with variational extensions can handle missing modalities gracefully
- **For Interpretability:** Cross-modal attention provides insights into which features the model focuses on

**Future Directions:** Modern systems often combine multiple strategies hierarchically. For example, using contrastive pre-training to learn initial representations, then fine-tuning with cross-modal attention for specific tasks, or employing autoencoder fusion for robustness to missing data while using attention mechanisms for interpretability.

## **Deep Learning Fusion Strategies - Comprehensive Guide**

Understanding and implementing multi-modal deep learning  
approaches

# Attention-based Integration



## Self-attention Fusion

Learning importance weights within modalities

## Cross-attention

Attention between different data modalities

## Multi-head Attention

Multiple attention perspectives simultaneously

## Hierarchical Attention

Feature and sample-level attention

## Interpretability

Attention weights as biological insights

# Detailed Attention Mechanisms

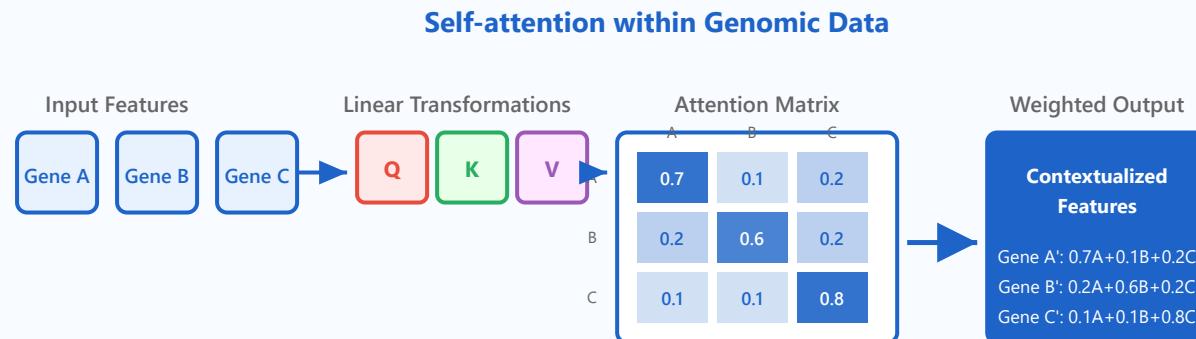
## 1. Self-attention Fusion

**Overview:** Self-attention fusion enables the model to learn dynamic importance weights for features within a single modality. Unlike fixed weighting schemes, self-attention adaptively determines which features are most relevant for the prediction task at hand.

**Core Mechanism:** The self-attention mechanism computes relationships between all positions in a sequence or feature set. For a given input, each feature acts as a query, key, and value simultaneously, allowing the model to weigh the relative importance of each feature based on its context.

$$\text{Self-Attention}(X) = \text{softmax}(QK^T / \sqrt{d_k})V$$

where  $Q = XW_Q$ ,  $K = XW_K$ ,  $V = XW_V$



### Biological Interpretation

- Gene A focuses mainly on itself (0.7) → Independent marker
- Gene B considers both A and itself → Co-regulatory relationship
  - Gene C is highly self-focused (0.8) → Distinct pathway

### Clinical Example: Gene Expression Analysis

In cancer genomics, self-attention can identify co-expressed gene modules. For instance, when analyzing breast cancer subtypes, the model might learn that BRCA1, BRCA2, and related DNA repair genes attend strongly to each other, revealing their co-regulatory network without prior biological knowledge.

### **Key Advantages:**

- **Adaptive weighting:** Importance scores change based on context and sample characteristics
- **Long-range dependencies:** Can capture relationships between distant features in the sequence
- **Parallel computation:** All attention scores computed simultaneously, enabling efficient processing
- **Biological discovery:** Learned attention patterns can reveal novel feature interactions

## **2. Cross-attention**

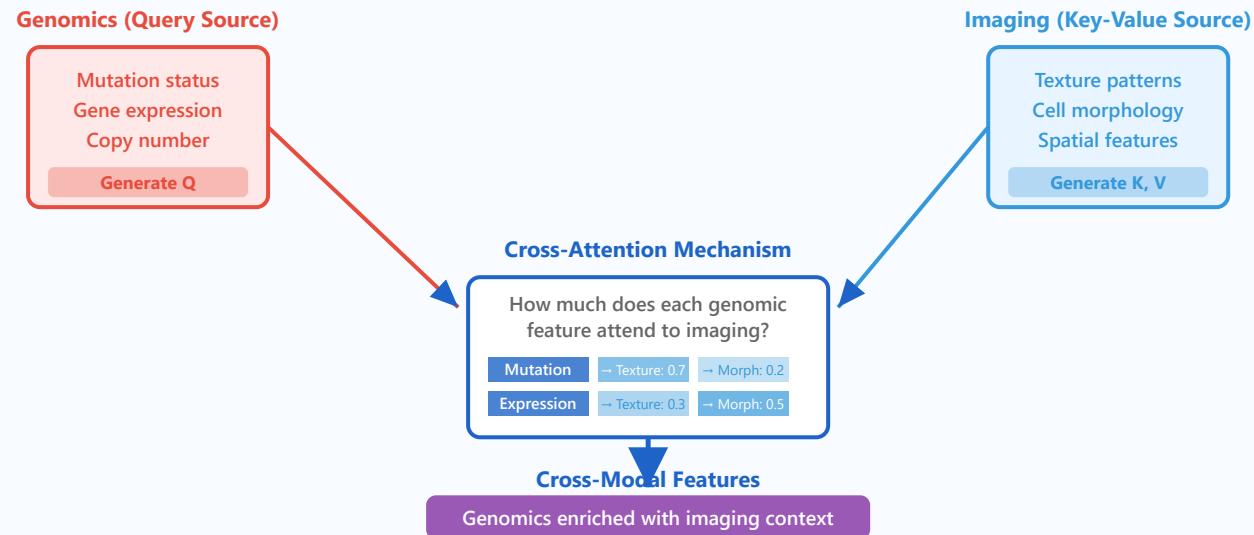
**Overview:** Cross-attention enables interaction between different data modalities by allowing one modality to attend to another. This mechanism is crucial for multi-modal learning where relationships between heterogeneous data sources must be discovered and leveraged.

**Core Mechanism:** Unlike self-attention where queries, keys, and values come from the same source, cross-attention uses one modality as the query source and another as the key-value source. This asymmetric attention allows directed information flow between modalities.

$$\text{Cross-Attention}(X, Y) = \text{softmax}(Q_X K_Y^T / \sqrt{d_k})V_Y$$

where Q comes from modality X, K and V come from modality Y

## Cross-attention: Imaging ← Genomics



### Clinical Example: Radiology-Pathology Integration

In lung cancer diagnosis, cross-attention can link CT imaging features with genomic alterations. The model might learn that certain texture patterns in CT scans (query from imaging) correlate strongly with EGFR mutation status (attending to genomic data), enabling non-invasive mutation prediction from imaging alone.

#### Key Features:

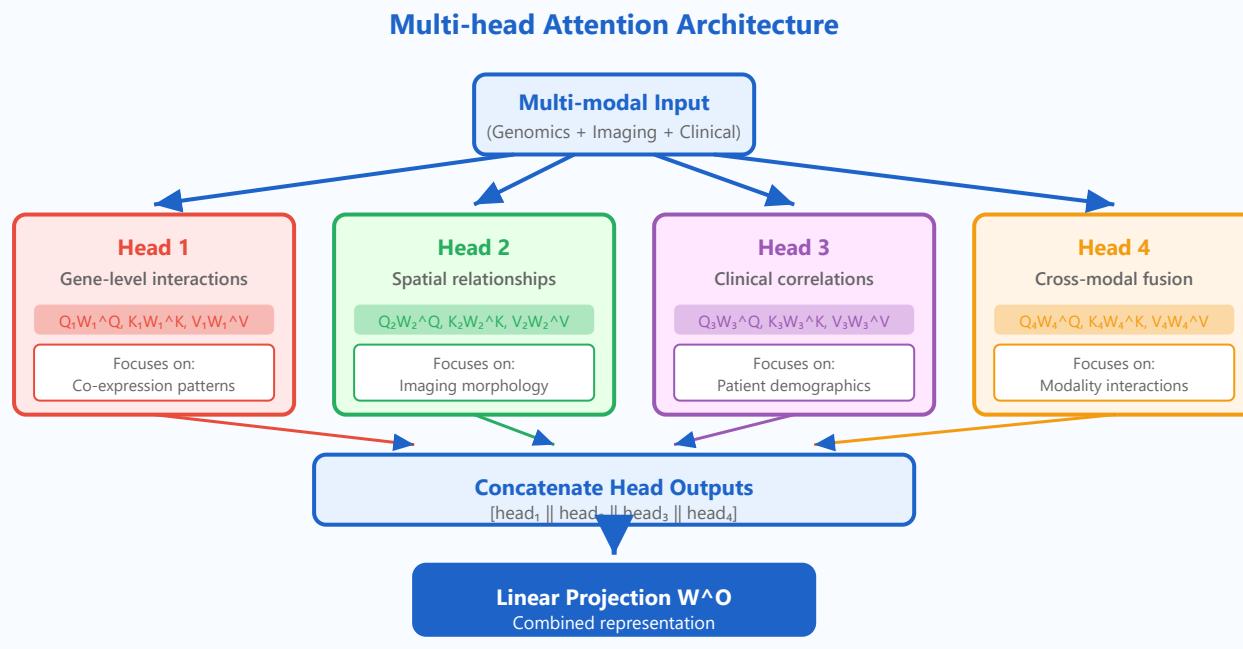
- **Bidirectional learning:** Can implement attention in both directions ( $A \rightarrow B$  and  $B \rightarrow A$ ) for symmetric relationships
- **Modality alignment:** Learns implicit alignment between heterogeneous feature spaces
- **Missing data handling:** One modality can query another even when partially observed
- **Clinical utility:** Enables prediction of expensive tests from cheaper modalities

### 3. Multi-head Attention

**Overview:** Multi-head attention runs multiple attention mechanisms in parallel, each learning different aspects of the relationships in the data. This is analogous to having multiple expert reviewers examine the same data from different perspectives.

**Core Mechanism:** Instead of performing a single attention function, multi-head attention linearly projects the queries, keys, and values  $h$  times with different learned projections. Each projection (head) captures different types of relationships, and their outputs are concatenated and linearly transformed.

```
MultiHead(Q, K, V) = Concat(head1, ..., headh)WO
where headi = Attention(QWi^Q, KWi^K, VWi^V)
```



#### Clinical Example: Cancer Subtype Classification

When classifying cancer subtypes from multi-modal data, different attention heads might specialize: Head 1 focuses on immune-related genes, Head 2 identifies tumor morphology patterns, Head 3 examines treatment history, and Head 4 integrates these perspectives. This parallel specialization captures the multi-faceted nature of cancer biology.

#### Benefits of Multiple Heads:

- **Diverse perspectives:** Each head can specialize in different relationship types or feature subsets
- **Robustness:** Multiple representations reduce reliance on any single attention pattern
- **Richer features:** Concatenated outputs provide more comprehensive representations
- **Biological alignment:** Different heads often correspond to distinct biological mechanisms

## 4. Hierarchical Attention

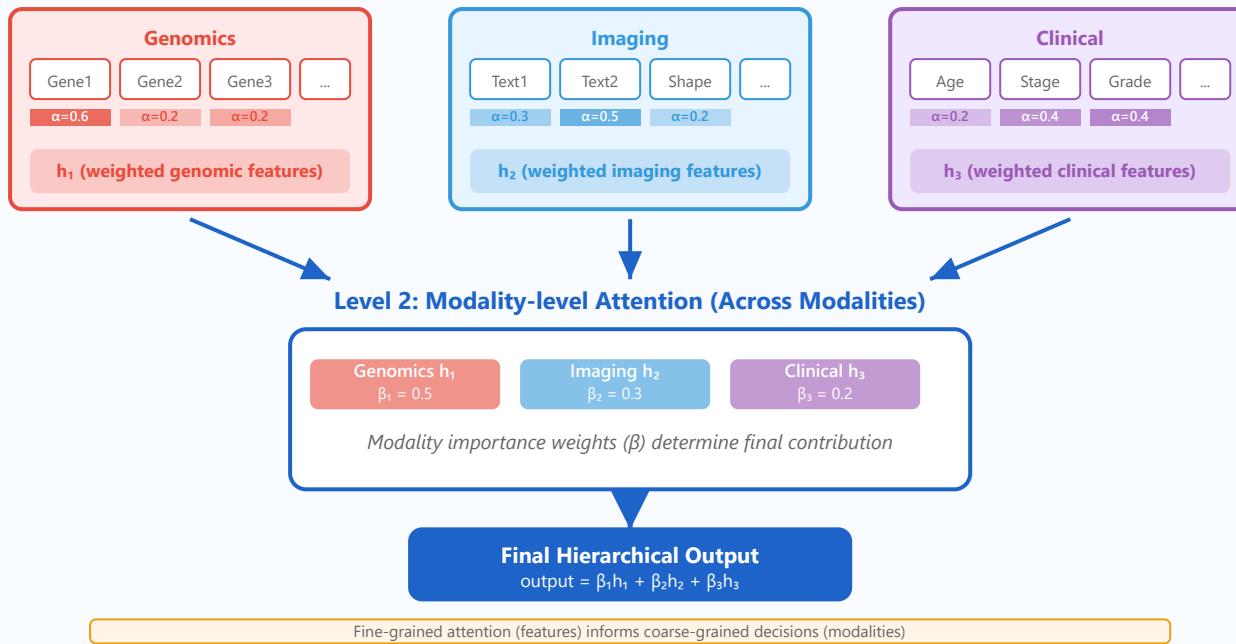
**Overview:** Hierarchical attention implements attention at multiple levels of granularity, typically combining feature-level attention with sample-level attention. This multi-scale approach mirrors the hierarchical structure of biological systems and clinical decision-making.

**Core Mechanism:** The hierarchical structure first applies attention at the feature level within each modality, then applies attention at the sample or modality level. This creates a two-stage (or multi-stage) process where fine-grained patterns inform coarse-grained decisions.

Level 1 (Feature):  $h_i = \sum_j \alpha_{ij} \cdot f_j$  (within modality i)  
Level 2 (Modality):  $output = \sum_i \beta_i \cdot h_i$  (across modalities)

## Hierarchical Attention: Two-Level Structure

### Level 1: Feature-level Attention (Within Each Modality)



### Clinical Example: Treatment Response Prediction

For predicting chemotherapy response, Level 1 attention identifies important genes within genomic data (e.g., DNA repair genes) and critical imaging features (e.g., tumor vascularity). Level 2 attention then determines that genomic features should receive higher weight ( $\beta_1=0.6$ ) than imaging ( $\beta_2=0.3$ ) for this particular patient, based on their molecular profile.

### Hierarchical Advantages:

- **Multi-scale reasoning:** Captures both fine-grained and coarse-grained patterns in a unified framework
- **Interpretability:** Two-level structure makes it easier to understand which features and modalities drive predictions
- **Efficiency:** Can prune less important features early in the hierarchy, reducing computation

- **Biological plausibility:** Mirrors hierarchical organization of biological systems and medical reasoning

## 5. Interpretability

---

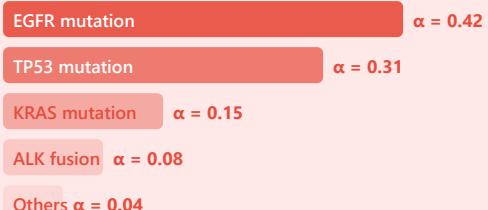
**Overview:** One of the most valuable aspects of attention mechanisms is their inherent interpretability. Attention weights provide direct insight into which features, samples, or modalities the model considers important for its predictions, making the decision-making process more transparent.

**Clinical Relevance:** In healthcare applications, model interpretability is not just desirable—it's essential. Clinicians need to understand why a model makes specific predictions to trust its recommendations, validate its reasoning against medical knowledge, and identify potential biases or errors.

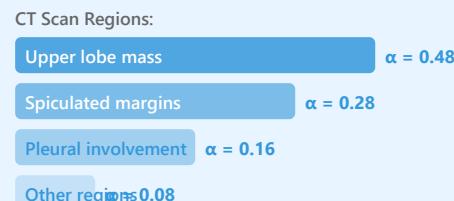
## Attention-based Interpretability in Cancer Diagnosis

Case: 62-year-old patient with suspected lung cancer → Model prediction: High-risk adenocarcinoma (95% confidence)

### Genomic Features Attention



### Imaging Features Attention



### Modality-level Attention Weights



### Clinical Interpretation & Validation

- EGFR mutation (42% attention) → Aligns with treatment guidelines for targeted therapy
- Upper lobe mass (48% imaging attention) → Confirms radiologist's primary concern
- Genomics prioritized (50%) → Appropriate for molecular subtyping in adenocarcinoma

### Real-world Application: Model Debugging

A model for predicting sepsis risk showed high attention to "patient room number" in one deployment. This attention weight visualization immediately revealed a data leakage problem—ICU patients (higher risk) had room numbers in a specific range. Without attention interpretability, this spurious correlation might have gone undetected, leading to poor generalization.

### Interpretability Benefits:

- **Model validation:** Verify that important features align with established medical knowledge
- **Error diagnosis:** Identify when the model focuses on spurious correlations or irrelevant features
- **Trust building:** Help clinicians understand and trust model predictions through transparent reasoning
- **Knowledge discovery:** Reveal novel feature combinations or interactions not previously recognized
- **Regulatory compliance:** Support model explainability requirements in healthcare AI regulations

- **Personalized insights:** Show patient-specific factors driving individual predictions

**Best Practice:**

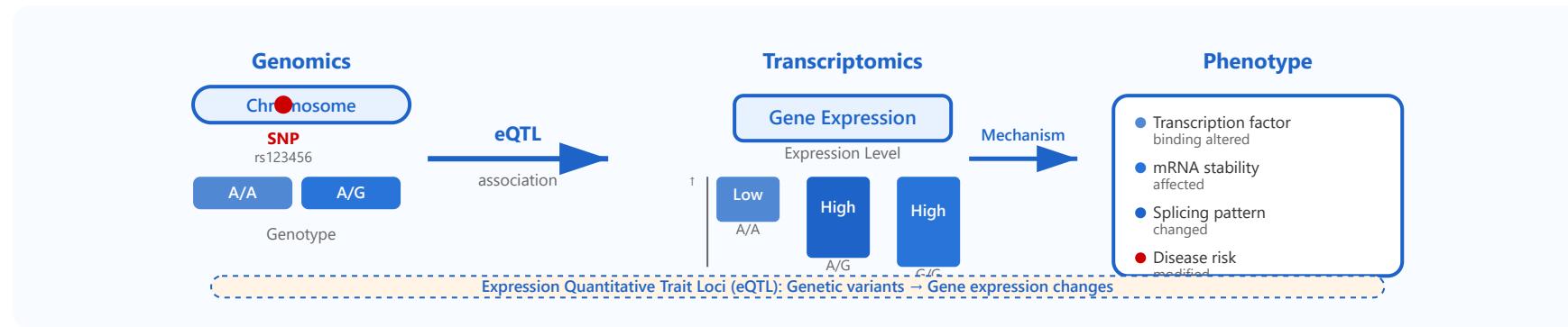
Always visualize attention weights during model development AND deployment  
Compare attention patterns across patient subgroups to identify biases

**Part 2/3:**

## **Multi-Omics Applications**

- Biological integration
- Technical challenges
- Analysis workflows

# Genomics + Transcriptomics Integration



## eQTL Analysis

Expression quantitative trait loci mapping

## ASE Detection

Allele-specific expression patterns

## Splicing QTLs

Genetic variants affecting RNA splicing

## Regulatory Variants

Non-coding variants and gene expression

## Allele-specific Binding

Transcription factor binding affected by SNPs

## Detailed Analysis of Integration Methods

## 1. eQTL Analysis (Expression Quantitative Trait Loci)

**Definition:** eQTLs are genomic loci that contribute to variation in gene expression levels. They represent genetic variants (typically SNPs) that are associated with changes in mRNA expression.

### Types of eQTLs:

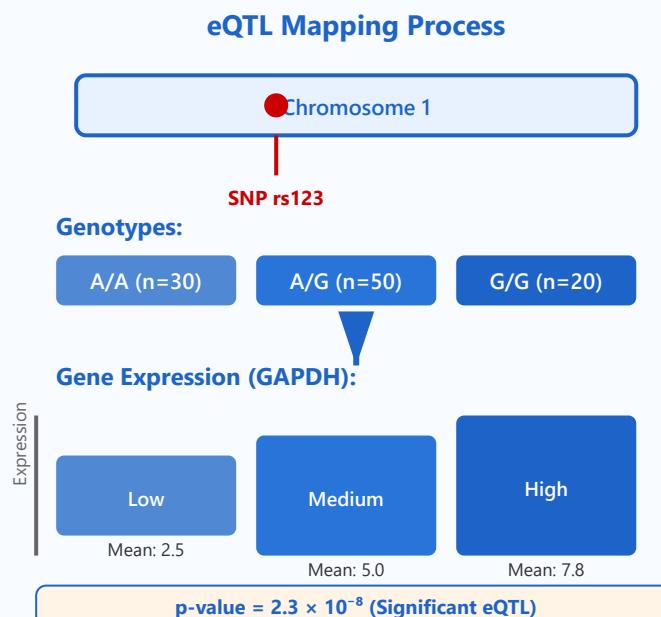
- **cis-eQTLs:** Variants located near the gene they regulate (typically within 1 Mb), acting locally on nearby genes
- **trans-eQTLs:** Variants located far from or on different chromosomes than the genes they regulate, acting through diffusible factors

### Applications:

- Identifying regulatory mechanisms underlying GWAS signals
- Prioritizing causal genes at disease-associated loci
- Understanding tissue-specific regulation
- Drug target identification and validation

### Statistical Approach:

Linear regression testing association between genotype and expression: Expression ~ Genotype + Covariates



### Key Points:

- eQTL studies require both genotype and gene expression data from the same individuals
- Most eQTLs are cis-acting and tissue-specific
- eQTL data helps interpret GWAS findings by linking variants to gene regulation
- Multiple testing correction is essential due to testing millions of SNP-gene pairs

## 2. ASE Detection (Allele-Specific Expression)

**Definition:** ASE occurs when the two alleles of a gene are expressed at different levels in heterozygous individuals. This reveals cis-regulatory effects at the individual level.

### Mechanisms Causing ASE:

- **Regulatory variants:** SNPs in promoters, enhancers, or UTRs affecting transcription
- **Genomic imprinting:** Parent-of-origin-specific expression
- **X-chromosome inactivation:** Random silencing of one X chromosome in females
- **Somatic mutations:** Acquired changes affecting expression

### Detection Methods:

- RNA-seq with allelic read counting at heterozygous SNPs

- Statistical testing (binomial or beta-binomial models)
- Requires phasing to determine which allele is expressed more

### Clinical Significance:

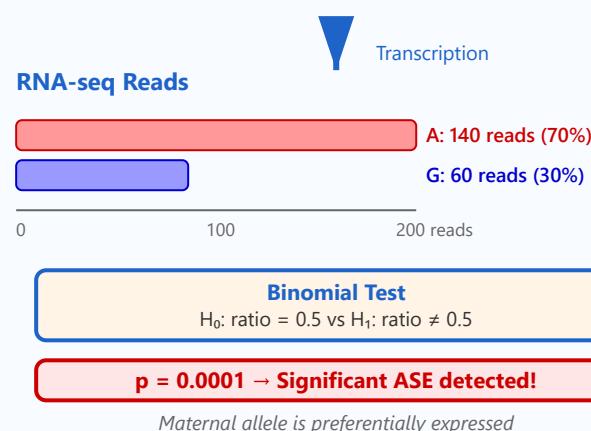
ASE can explain incomplete penetrance, variable expressivity, and personalized drug responses in genetic diseases.

### Allele-Specific Expression

#### Genomic DNA (Heterozygous)



Expected:  
50% : 50%  
**Observed:**  
70% : 30%



#### Key Points:

- ASE provides direct evidence of cis-regulatory variation
- Requires sufficient RNA-seq read depth at heterozygous sites (typically >20 reads)
- Can identify imprinted genes and parent-of-origin effects
- Important for understanding disease mechanisms and personalized medicine

## 3. Splicing QTLs (sQTLs)

**Definition:** sQTLs are genetic variants that affect RNA splicing patterns, leading to different isoforms or exon

usage patterns. They represent a major mechanism of gene regulation.

### Types of Splicing Changes:

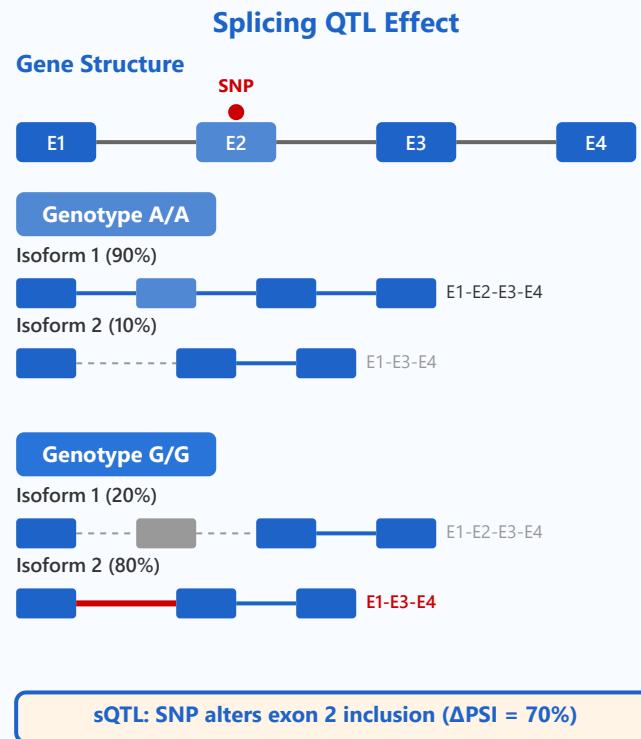
- **Exon skipping/inclusion:** Complete exons included or excluded
- **Alternative 5' or 3' splice sites:** Different splice site usage
- **Intron retention:** Failure to remove introns
- **Alternative first/last exons:** Different transcription start/end sites

### Molecular Mechanisms:

- Disruption of splice site consensus sequences
- Creation or destruction of splicing enhancers/silencers
- Alteration of RNA secondary structure
- Changes in splicing factor binding sites

### Detection Methods:

LeafCutter, MAJIQ, or rMATS for differential splicing analysis; testing association between genotype and junction usage or PSI (Percent Spliced In) values.



### Key Points:

- sQTLs are highly prevalent and affect ~30% of genes
- Often independent of expression-level effects (eQTLs)
- Can create functionally distinct protein isoforms

- Important for understanding disease mechanisms and therapeutic targets
- Tissue-specific splicing patterns reflect cell-type-specific regulatory programs

## 4. Regulatory Variants in Non-Coding Regions

**Definition:** Regulatory variants are genetic variations in non-coding regions that affect gene expression by altering transcriptional regulation. Most disease-associated variants identified by GWAS are in non-coding regions.

### Types of Regulatory Elements:

- **Promoters:** Regions directly upstream of genes controlling transcription initiation
- **Enhancers:** Distal elements that increase transcription (can be >1 Mb away)
- **Silencers:** Elements that repress transcription
- **Insulators:** Boundary elements that prevent enhancer-promoter interactions
- **3' UTRs:** Regions affecting mRNA stability and translation

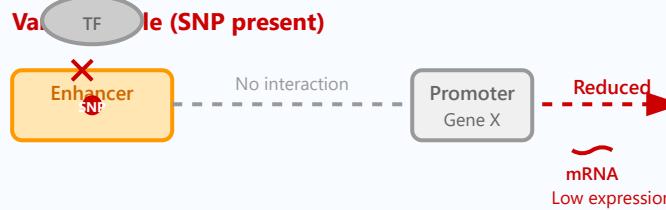
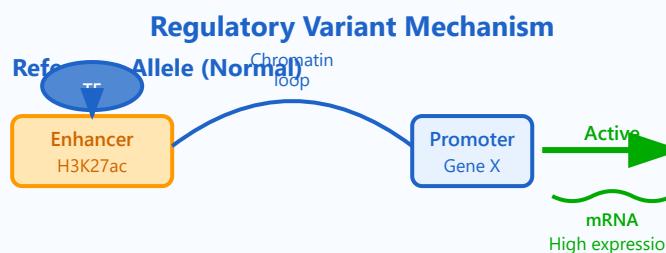
### Identification Strategies:

- Integration with ChIP-seq data (histone marks, TF binding)

- ATAC-seq or DNase-seq for chromatin accessibility
- Chromosome conformation capture (Hi-C, ChIA-PET) for enhancer-promoter interactions
- Massively parallel reporter assays (MPRA)

### Functional Consequences:

Altered transcription factor binding, changes in chromatin accessibility, disrupted enhancer-promoter loops, modified histone modifications.



**SNP disrupts TF binding site → Reduced enhancer activity**  
→ Decreased gene expression → Disease phenotype

### Key Points:

- ~90% of disease-associated variants from GWAS are in non-coding regions
- Enhancers can regulate genes over long distances (>1 Mb)
- Regulatory effects are often cell-type and tissue-specific
- Integration with epigenomic data is essential for interpretation
- Functional validation requires reporter assays or genome editing

## 5. Allele-Specific Transcription Factor Binding

**Definition:** Allele-specific binding (ASB) occurs when a transcription factor (TF) preferentially binds to one allele over another at heterozygous SNP sites. This is a direct mechanism linking genetic variation to gene regulation.

### Detection Methods:

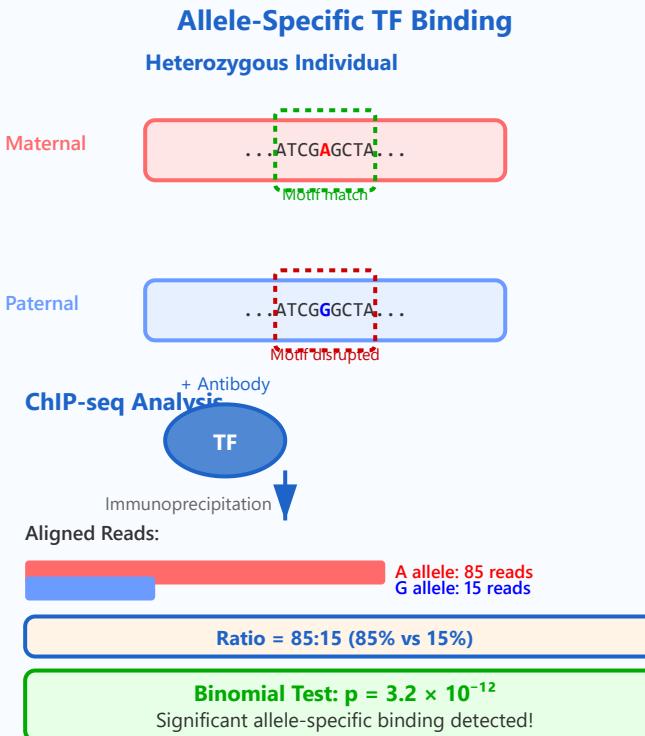
- **ChIP-seq:** Chromatin immunoprecipitation followed by sequencing to map TF binding
- **Analysis:** Count reads mapping to each allele at heterozygous SNPs within ChIP-seq peaks
- **Statistical test:** Binomial test or beta-binomial model for differential binding

### Mechanisms:

- Direct disruption of TF binding motif by SNP
- Indirect effects through DNA shape changes
- Altered cooperativity with other TFs
- Changes in local chromatin structure

### Biological Significance:

ASB events link genetic variants to gene expression changes and explain mechanisms of disease-associated variants. They provide functional evidence for causal variants in GWAS loci.



### Key Points:

- ASB directly demonstrates functional impact of genetic variants on TF binding
- Requires sufficient ChIP-seq read depth at heterozygous sites (>20 reads)

- Can identify causal variants at GWAS loci by showing mechanism
- Often correlates with allele-specific expression (ASE) at target genes
- Useful for fine-mapping disease-associated variants and drug target discovery

# Proteogenomics: Integrating Genomics and Proteomics

---

## Variant Peptides

Protein sequences from genomic variants

## Novel ORFs

Discovering new protein-coding regions

## PTM Sites

Post-translational modification mapping

## Protein Isoforms

Alternative splicing products in proteomics

## Neo-antigens

Tumor-specific antigens for immunotherapy

1

## Variant Peptides

**Reference DNA:**

ATG GCC TAT GGC AAA TTC GGA

SNP (T)  
**Variant DNA:**

ATG GCC C AT GGC AAA TTC GGA

**Reference Protein:**

Met-Ala-Tyr-Gly-Lys-Phe-Gly

**Variant Protein:**

Variant peptides represent protein sequences that arise from genomic variations such as single nucleotide polymorphisms (SNPs), insertions, deletions, or other mutations. Proteogenomics allows us to identify these variant peptides by integrating genomic sequencing data with mass spectrometry-based proteomics.

**How It Works:**

Genomic variants identified through DNA or RNA sequencing are used to create customized protein sequence databases. These databases include both reference sequences and variant sequences. Mass spectrometry data is then searched against these expanded databases to identify peptides that match variant sequences rather than the reference genome.

**Key Points:**

- SNPs can lead to amino acid substitutions, creating variant peptides
- Indels (insertions/deletions) may cause frameshift mutations
- Proteogenomics validates genomic predictions at the protein level
- Essential for personalized medicine and pharmacogenomics

## Clinical Applications:

- **Pharmacogenomics:** Understanding how genetic variants affect drug metabolism
- **Disease susceptibility:** Identifying protein variants associated with disease risk
- **Biomarker discovery:** Finding variant peptides as diagnostic markers

## 2 Novel Open Reading Frames (ORFs)

### Conventional Annotation:



### Ribosome Profiling + Mass Spec:



### Discoveries:

Upstream ORF (uORF)

Protein-coding lncRNA

Alternative start codon

Novel open reading frames (ORFs) are previously unannotated protein-coding sequences discovered through proteogenomic approaches. Traditional genome annotation relies heavily on computational predictions, which can miss short ORFs, alternative start codons, or protein-coding sequences in regions previously classified as non-coding.

## Discovery Methods:

Proteogenomics combines ribosome profiling (which shows where ribosomes are actively translating) with mass spectrometry to provide direct evidence of protein translation from novel ORFs. This approach has revealed that genomes are more complex than previously thought, with many small proteins and micropeptides being actively expressed.

### Types of Novel ORFs:

- **Upstream ORFs (uORFs):** Short coding sequences in the 5' UTR that regulate main ORF translation
- **Downstream ORFs (dORFs):** Additional coding sequences in the 3' UTR
- **Long non-coding RNA (lncRNA) ORFs:** Protein-coding capacity in presumed non-coding transcripts
- **Short ORFs (sORFs):** Micropeptides under 100 amino acids often missed by annotation pipelines
- **Alternative translation start sites:** Non-AUG start codons producing protein variants

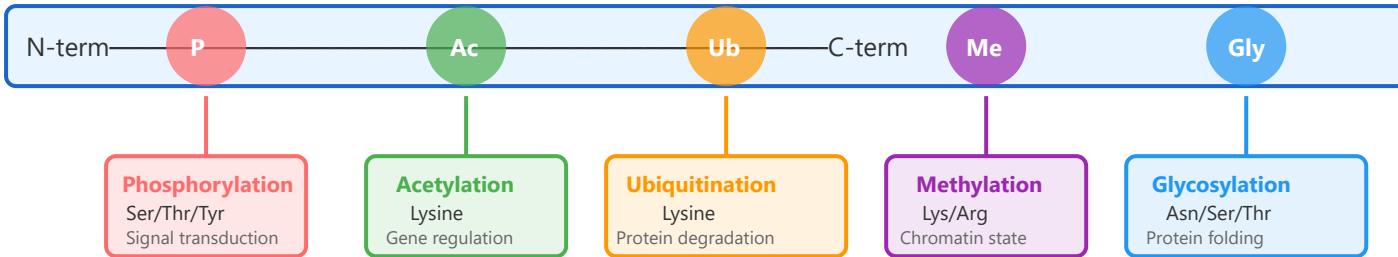
### Research Impact:

- **Genome reannotation:** Improving the accuracy of gene catalogs
- **Functional genomics:** Understanding previously unknown regulatory mechanisms
- **Drug target discovery:** Identifying new therapeutic targets among novel proteins
- **Evolution studies:** Understanding the birth and evolution of new genes

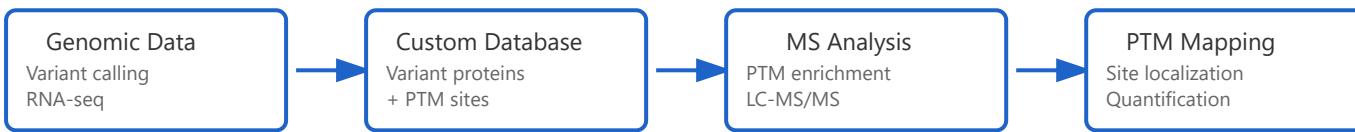
3

## Post-Translational Modification (PTM) Sites

### Protein with Multiple PTMs:



### Proteogenomic PTM Mapping:



Post-translational modifications (PTMs) are chemical modifications to proteins after translation that dramatically expand protein functional diversity. Proteogenomics enhances PTM analysis by providing genomic context, allowing researchers to map PTM sites to specific protein variants, isoforms, and novel ORFs.

### Major PTM Types and Functions:

Over 400 different types of PTMs have been identified, but the most commonly studied include phosphorylation, acetylation, methylation, ubiquitination, and glycosylation. Each PTM type has distinct regulatory functions and can profoundly affect protein activity, localization, stability, and interactions.

#### Proteogenomic Advantages for PTM Studies:

- Variant-specific PTMs:** Identify how genetic variants create or eliminate PTM sites
- Isoform-specific modifications:** Map PTMs to specific alternative splicing isoforms
- Novel site discovery:** Find PTMs in previously unannotated proteins

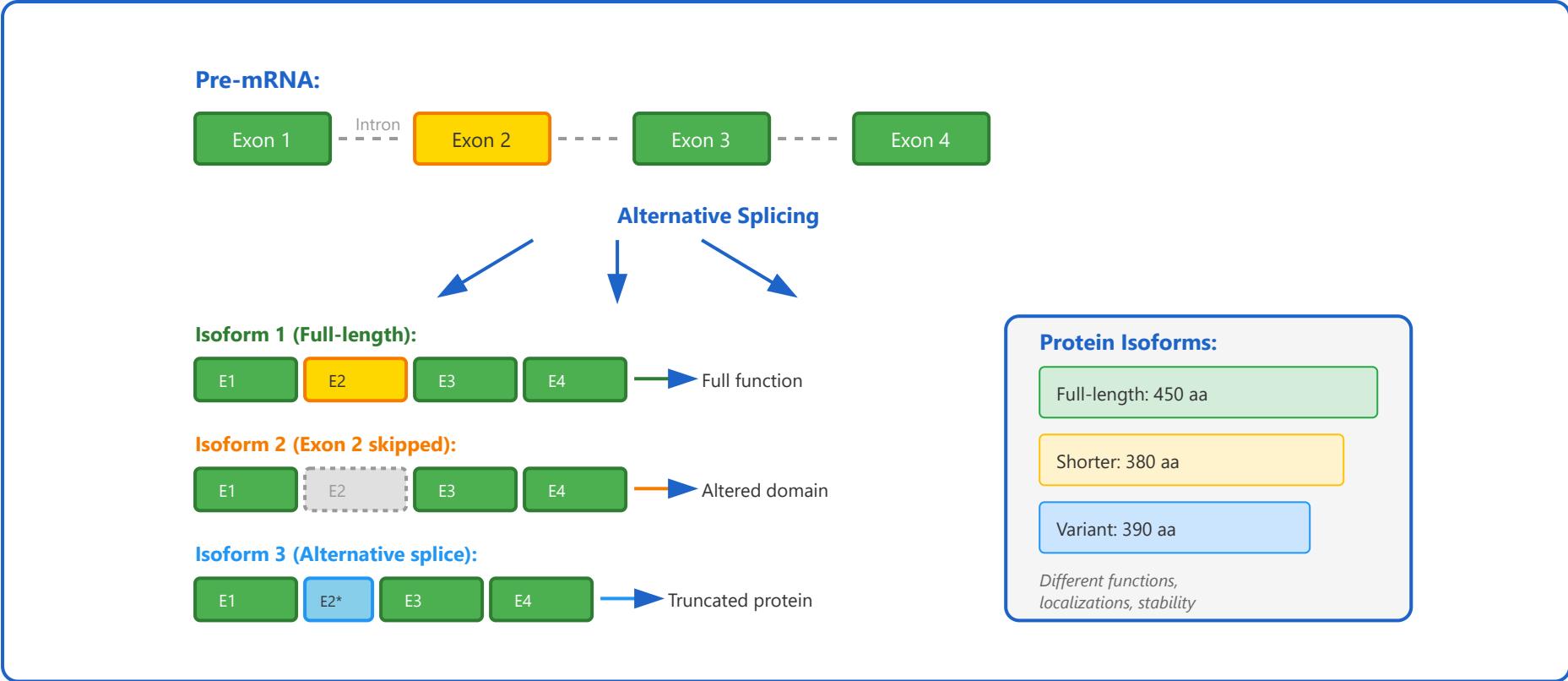
- **Functional context:** Link PTMs to genomic features and disease variants
- **Dynamic regulation:** Track PTM changes across conditions or disease states

#### Clinical and Research Applications:

- **Cancer biology:** Aberrant PTM patterns as hallmarks of cancer
- **Kinase inhibitor development:** Targeting phosphorylation-dependent signaling
- **Epigenetics:** Understanding histone modifications and gene regulation
- **Neurodegenerative diseases:** Role of aberrant ubiquitination and phosphorylation
- **Biomarker discovery:** PTM signatures for disease diagnosis and prognosis

4

## Protein Isoforms from Alternative Splicing



Alternative splicing is a fundamental mechanism that generates multiple protein isoforms from a single gene. It's estimated that over 95% of human multi-exon genes undergo alternative splicing, dramatically expanding proteomic diversity. Proteogenomics is essential for validating these isoforms at the protein level and understanding their functional consequences.

### Types of Alternative Splicing:

Alternative splicing can occur through several mechanisms: exon skipping (cassette exons), alternative 5' or 3' splice sites, intron retention, mutually exclusive exons, and alternative promoters or polyadenylation sites. Each mechanism produces distinct protein isoforms with potentially different functions, cellular localizations, or regulatory properties.

### Proteogenomic Challenges and Solutions:

- **Isoform inference:** RNA-seq identifies splice junctions; MS validates protein expression
- **Isoform-specific peptides:** Unique peptides that distinguish between isoforms

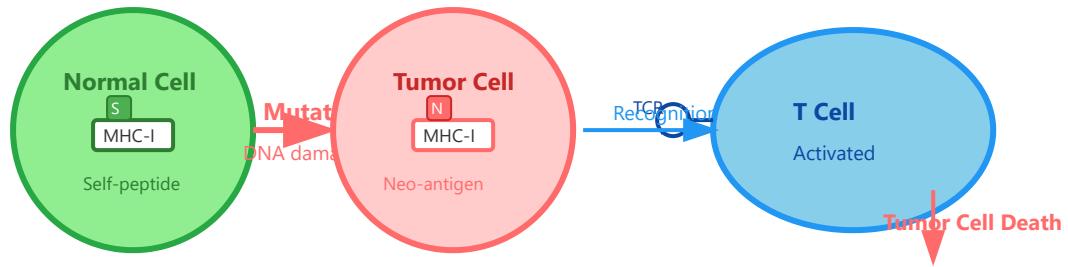
- **Quantification:** Measuring relative abundance of different isoforms
- **Functional annotation:** Linking isoforms to specific biological functions
- **Disease relevance:** Aberrant splicing in cancer and genetic diseases

#### Biological and Clinical Significance:

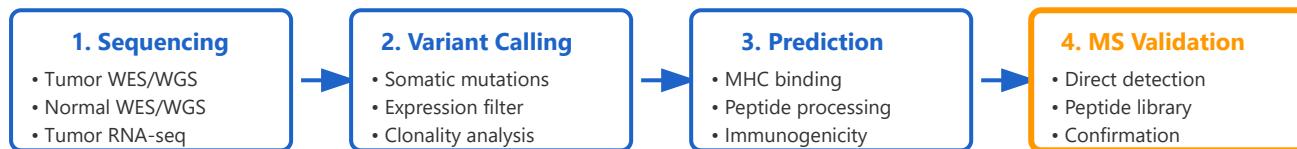
- **Tissue-specific expression:** Different isoforms predominate in different tissues
- **Development and differentiation:** Isoform switching during cell fate transitions
- **Cancer diagnostics:** Aberrant isoform ratios as cancer biomarkers
- **Therapeutic targeting:** Isoform-specific drugs for precision medicine
- **Splice-modulating therapies:** ASOs and small molecules to correct splicing defects

5

## Neo-antigens for Cancer Immunotherapy



### Neo-antigen Discovery Pipeline:



#### Therapeutic Applications

Personalized Vaccines

Adoptive T Cell Therapy

Checkpoint Blockade + Neo

Neo-antigens are tumor-specific antigens arising from somatic mutations in cancer cells. These novel peptide sequences are not present in normal tissues, making them ideal targets for cancer immunotherapy. Proteogenomics plays a crucial role in neo-antigen discovery by combining genomic mutation data with proteomic validation to identify peptides that are actually presented on tumor cell surfaces.

### The Neo-antigen Pipeline:

Neo-antigen discovery begins with identifying somatic mutations through whole-exome or whole-genome sequencing of tumor and matched normal tissue. RNA-seq data filters for expressed mutations. Computational tools predict which mutant peptides will bind to the patient's specific MHC molecules and be presented on the cell surface. Critically, mass spectrometry provides direct experimental evidence that predicted neo-antigens are actually present on tumor cells.

### **Types of Neo-antigens:**

- **SNV-derived:** Point mutations creating single amino acid substitutions
- **Indel-derived:** Frameshift mutations generating completely novel sequences
- **Gene fusion neo-antigens:** Peptides spanning fusion breakpoints
- **Splicing-derived:** Aberrant splicing creating novel junctions
- **Non-coding neo-antigens:** Peptides from presumed non-coding regions

### **Why Mass Spectrometry Validation Matters:**

While computational prediction identifies thousands of potential neo-antigens, only a small fraction are actually processed, presented on MHC molecules, and detectable by mass spectrometry. MS validation reduces false positives, confirms peptide processing, validates MHC presentation, and prioritizes the most promising candidates for therapeutic development.

### **Clinical Applications and Impact:**

- **Personalized cancer vaccines:** Patient-specific vaccines targeting validated neo-antigens
- **CAR-T and TCR-T therapy:** Engineering T cells to recognize specific neo-antigens
- **Checkpoint inhibitor biomarkers:** Tumor mutational burden and neo-antigen load predict response
- **Combination strategies:** Neo-antigen vaccines combined with checkpoint blockade
- **Minimal residual disease monitoring:** Tracking neo-antigen-specific immune responses

### **Success Stories and Future Directions:**

- Clinical trials showing personalized neo-antigen vaccines can induce specific T cell responses
- Improved response rates when combined with checkpoint inhibitors
- Development of off-the-shelf shared neo-antigen libraries for common cancers
- Integration of AI/ML for better neo-antigen prediction and prioritization

- Expanding beyond MHC-I to MHC-II neo-antigens for CD4+ T cell activation

# Imaging-genomics (Radiogenomics)

## Radiogenomics

Linking imaging phenotypes to genotypes

## Imaging Features

Quantitative features from medical images

## Genetic Associations

GWAS-style analysis with imaging

## Outcome Prediction

Combining imaging and genomics for prognosis

## Treatment Response

Predicting therapy efficacy

## • 1. Radiogenomics: Linking Imaging Phenotypes to Genotypes

**Radiogenomics** is an emerging field that integrates medical imaging data with genomic information to understand the relationship between imaging characteristics and underlying genetic profiles.

### Core Concept

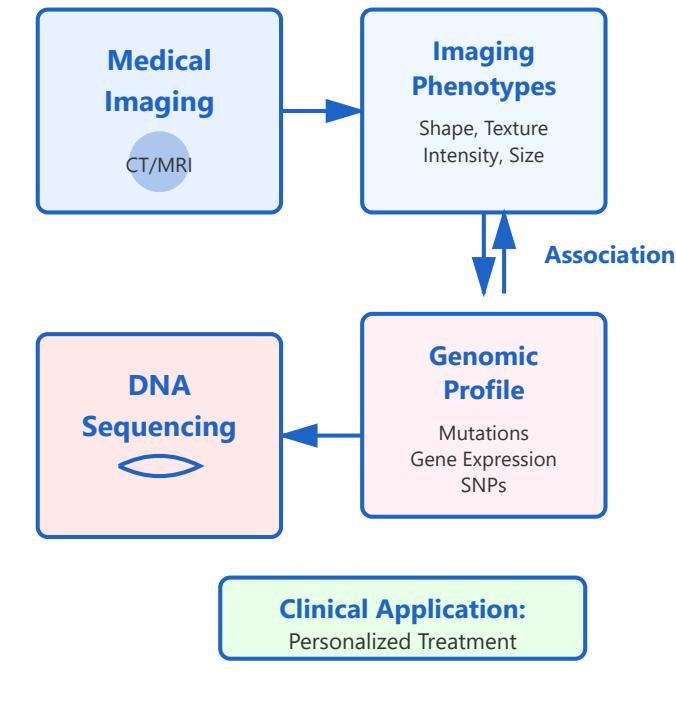
Radiogenomics explores how genetic variations influence imaging phenotypes (observable characteristics in medical images) and how imaging features can reflect molecular and genetic alterations in tissues.

## Key Applications

- **Cancer characterization:** Identifying tumor subtypes based on imaging patterns linked to specific mutations (e.g., EGFR, KRAS in lung cancer)
- **Non-invasive genotyping:** Predicting genetic status without requiring tissue biopsy
- **Personalized medicine:** Tailoring treatment strategies based on imaging-genomic profiles
- **Disease mechanisms:** Understanding how genetic variations manifest as observable imaging characteristics

## Key Insight

Radiogenomics bridges the gap between macroscopic imaging and microscopic molecular biology, enabling non-invasive characterization of disease at the genetic level.

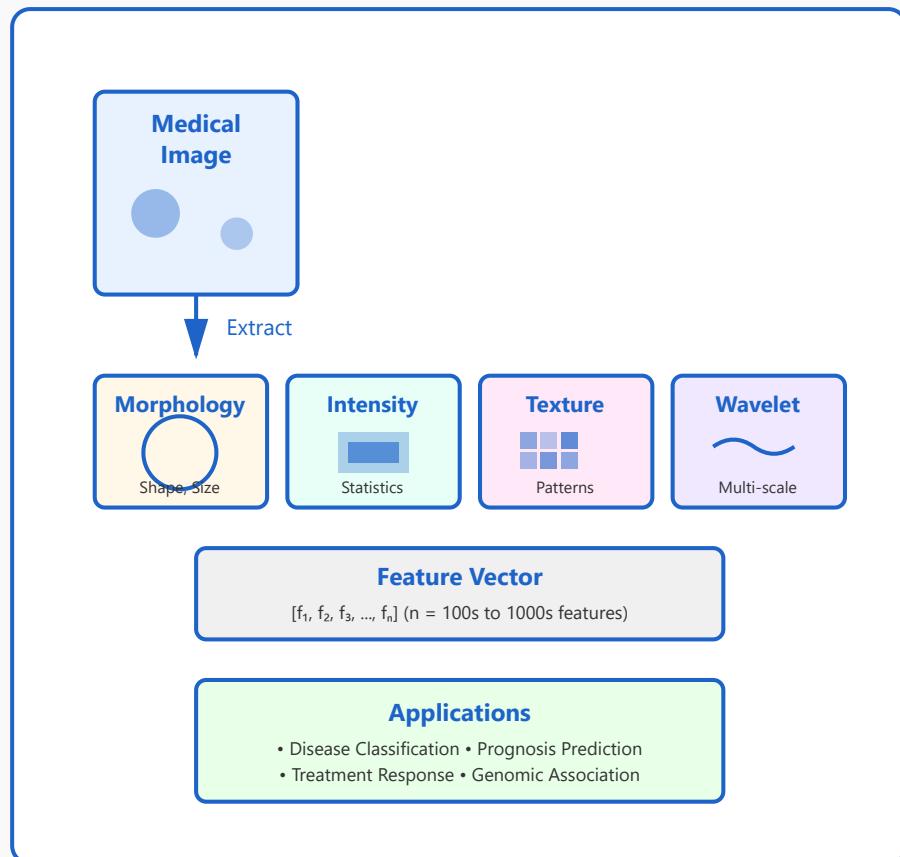


## • 2. Imaging Features: Quantitative Analysis of Medical Images

**Imaging features** are quantitative measurements extracted from medical images that characterize tissue properties, lesion characteristics, and spatial patterns.

### Types of Imaging Features

- **Morphological features:** Shape, size, volume, surface area, compactness, sphericity
- **Intensity-based features:** Mean, median, standard deviation, skewness, kurtosis
- **Textural features:** Gray-level co-occurrence matrix (GLCM), gray-level run-length matrix (GLRLM), entropy, homogeneity
- **Wavelet features:** Multi-scale decomposition of image signals
- **Functional features:** Perfusion, diffusion, metabolism from functional imaging (fMRI, PET)



### Radiomics

Radiomics is the high-throughput extraction of large numbers of quantitative features from medical images, transforming images into mineable data for clinical decision support.

### Clinical Value

These features can capture subtle patterns invisible to the human eye and provide objective, reproducible measurements for disease characterization and monitoring.

### • 3. Genetic Associations: GWAS-style Analysis with Imaging

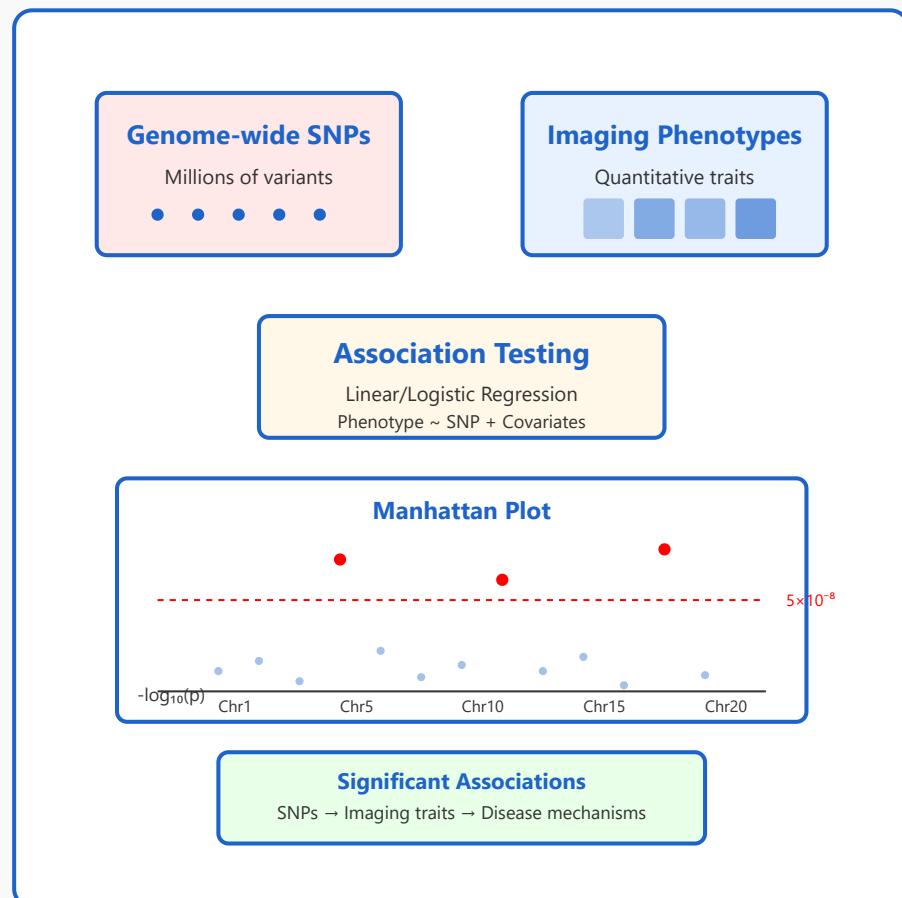
**Imaging genetics** combines genome-wide association study (GWAS) methodology with quantitative imaging phenotypes to identify genetic variants associated with brain structure, tumor characteristics, or other imaging features.

#### Methodology

Similar to traditional GWAS, imaging genetics studies test associations between millions of genetic variants (SNPs) and imaging-derived phenotypes across large cohorts.

#### Key Applications

- **Brain imaging genetics:** Identifying genes affecting brain volume, cortical thickness, white matter integrity (e.g., ENIGMA consortium studies)
- **Cancer radiogenomics:** Linking tumor imaging features to driver mutations and molecular subtypes



- **Cardiovascular imaging:** Genetic determinants of cardiac structure and function
- **Disease risk prediction:** Using genetic variants associated with imaging markers to predict disease susceptibility

## Statistical Considerations

Multiple testing correction (e.g., Bonferroni, FDR) is essential due to testing millions of SNPs. Typical significance threshold:  $p < 5 \times 10^{-8}$ .

### Impact

This approach has revealed genetic architecture of brain structure and identified novel disease mechanisms by connecting genetic variants to intermediate imaging phenotypes.

## • 4. Outcome Prediction: Integrating Imaging and Genomics for Prognosis

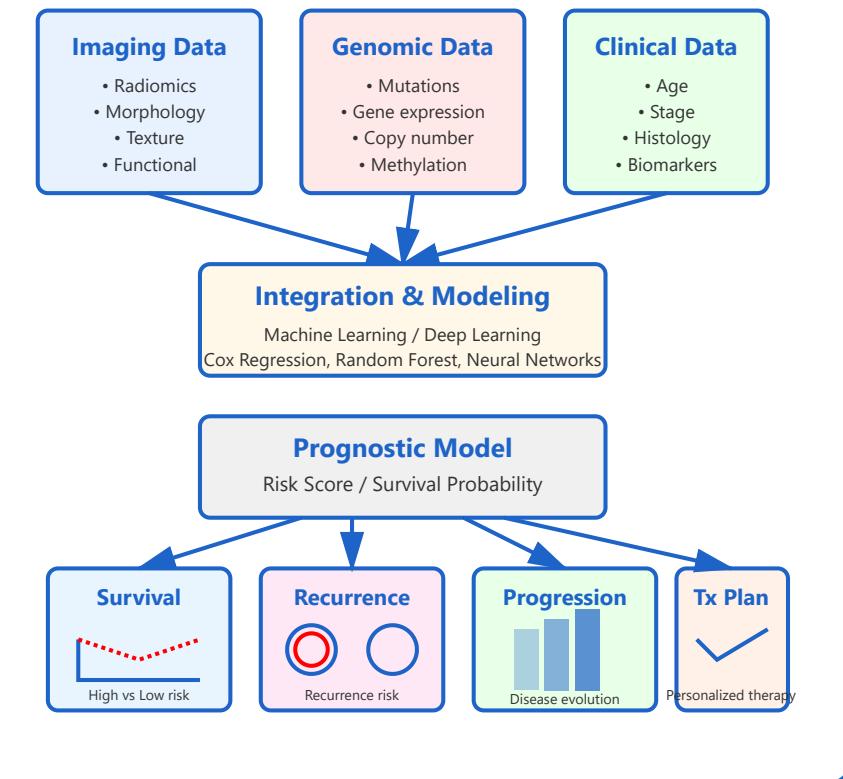
**Outcome prediction** combines imaging features and genomic data to forecast disease progression, survival, recurrence risk, and other clinical endpoints.

### Integration Approaches

- **Multi-modal models:** Machine learning models that integrate radiomics features with genomic profiles
- **Complementary information:** Imaging captures spatial heterogeneity; genomics reveals molecular drivers
- **Deep learning:** End-to-end models learning from both image data and genomic sequences

## Clinical Applications

- **Cancer prognosis:** Predicting overall survival, disease-free survival, metastasis risk
- **Risk stratification:** Identifying high-risk vs. low-risk patient groups
- **Recurrence prediction:** Forecasting likelihood of disease recurrence after treatment
- **Progression monitoring:** Tracking disease evolution over time



## Performance Metrics

Models are evaluated using concordance index (C-index), area under ROC curve (AUC), calibration plots, and decision curve analysis.

### Clinical Impact

Combined imaging-genomic models often outperform single-modality approaches, providing more accurate and

## • 5. Treatment Response: Predicting Therapy Efficacy

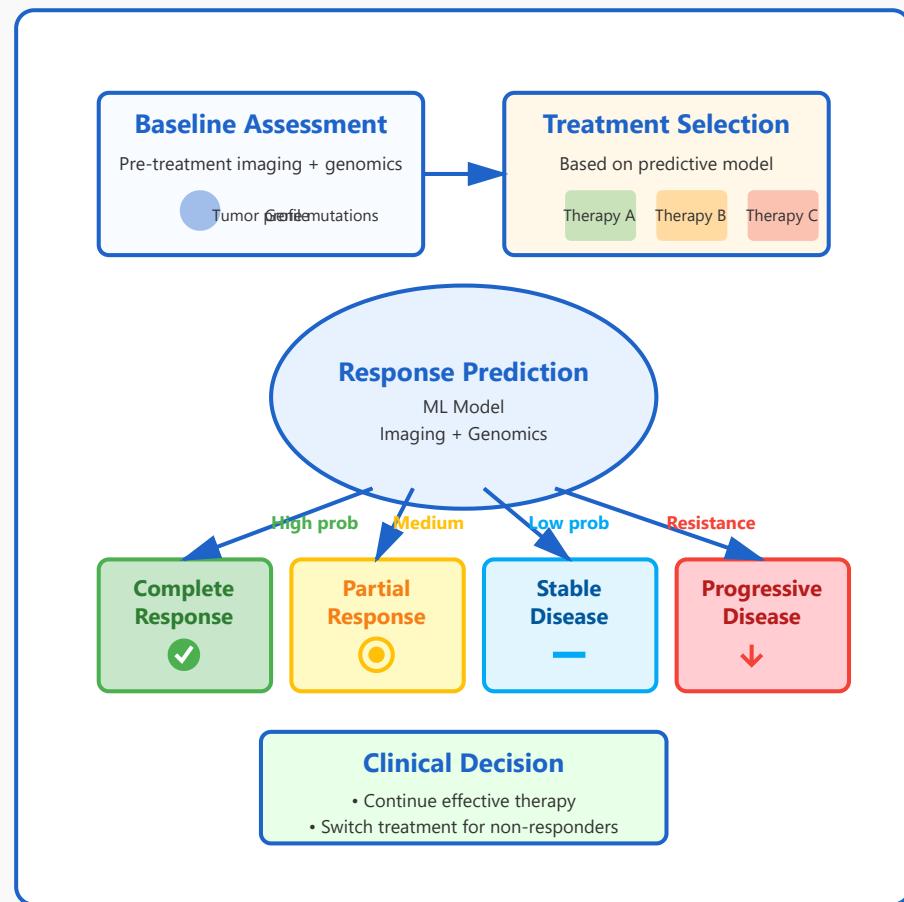
**Treatment response prediction** leverages baseline and longitudinal imaging-genomic data to forecast how patients will respond to specific therapies, enabling precision medicine approaches.

### Prediction Strategies

- **Baseline prediction:** Using pre-treatment imaging and genomic profiles to predict response
- **Early response assessment:** Analyzing early treatment changes to predict long-term outcomes
- **Resistance mechanisms:** Identifying imaging-genomic signatures of treatment resistance

### Clinical Applications

- **Chemotherapy response:** Predicting response to cytotoxic agents based on tumor heterogeneity and genomic profiles
- **Targeted therapy:** Matching patients to targeted agents (e.g., EGFR inhibitors for EGFR-mutant tumors with specific



imaging features)

- **Immunotherapy:** Identifying imaging biomarkers (e.g., tumor infiltration patterns) and genomic markers (e.g., PD-L1, TMB) predicting immunotherapy response
- **Radiation response:** Predicting radiosensitivity based on tumor oxygenation (imaging) and DNA repair genes

## Response Criteria

RECIST (Response Evaluation Criteria in Solid Tumors)  
combined with functional imaging changes and genomic evolution tracking.

### Personalized Treatment

This approach enables selection of optimal therapies for individual patients, avoiding ineffective treatments and their associated toxicities while improving outcomes.

# Clinical + Molecular Data Integration

## EHR Integration

Electronic health records with omics data

## Lab Values

Clinical laboratory measurements

## Imaging Reports

Radiology and pathology findings

## Molecular Profiles

Genomic, transcriptomic, proteomic data

## Temporal Alignment

Synchronizing time-series clinical and molecular data

1

## EHR Integration

Electronic Health Record (EHR) integration with omics data represents a paradigm shift in precision medicine, enabling comprehensive patient profiling that combines traditional clinical information with molecular-level insights.

### Key Components

- **Patient Demographics:** Age, sex, ethnicity, and family history integrated with genetic predisposition data

 Patient Demographics & History



 Genomic Data (WGS/WES)



 Gene Expression (RNA-seq)

- **Clinical History:** Past diagnoses, medications, and treatments correlated with molecular response patterns
- **Omics Layers:** Genomic variants, gene expression profiles, and epigenetic modifications
- **Standardization:** FHIR (Fast Healthcare Interoperability Resources) standards for data exchange

## Clinical Applications

- Pharmacogenomic decision support for drug selection and dosing
- Risk stratification based on genetic and clinical factors
- Early disease detection through molecular biomarkers
- Personalized treatment pathway recommendations

↓

Protein Data (Proteomics)

↓

Treatment Response Data

↓

✓ Integrated EHR-Omics Profile

## 2 Lab Values

Clinical laboratory measurements provide quantitative biomarkers that can be directly correlated with molecular profiles to reveal disease mechanisms and predict treatment outcomes.

### Laboratory Data Types

Hemoglobin  
12.5 g/dL

WBC  
8,200/ $\mu$ L

Glucose  
105 mg/dL

Creatinine  
1.1 mg/dL

ALT  
35 U/L

CRP  
2.5 mg/L

- **Hematology:** Complete blood count (CBC), differential counts, coagulation studies
- **Clinical Chemistry:** Metabolic panels, liver and kidney function tests, electrolytes
- **Immunology:** Cytokine levels, antibody titers, inflammatory markers
- **Tumor Markers:** PSA, CA-125, CEA, AFP for cancer monitoring
- **Molecular Tests:** PCR results, viral loads, ctDNA measurements

## Integration Strategies

- Correlate lab abnormalities with gene expression changes
- Identify molecular pathways underlying clinical phenotypes
- Monitor treatment response through serial measurements
- Predict lab value trends using molecular signatures

 PSA  
4.2 ng/mL

 ctDNA  
0.15%

Correlated with Gene Expression Profiles

## 3 Imaging Reports

Medical imaging provides structural and functional information that, when integrated with molecular data, enables radiogenomics—the study of relationships between imaging features and genomic patterns.

 Radiology Report: 3.2cm mass, irregular margins



## Imaging Modalities

- **CT/MRI:** Tumor size, morphology, and anatomical relationships
- **PET Scans:** Metabolic activity and molecular targeting (FDG-PET, PSMA-PET)
- **Pathology Imaging:** Digital histopathology and immunohistochemistry
- **Functional Imaging:** Perfusion, diffusion, and spectroscopy data

## Radiogenomics Applications

- Predict molecular subtypes from imaging features using AI
- Non-invasive assessment of tumor heterogeneity
- Correlate imaging biomarkers with gene signatures
- Monitor spatial-temporal evolution of disease
- Guide biopsy sites based on molecular likelihood

Pathology: Adenocarcinoma, Grade 2/3

IHC: ER+, PR+, HER2-

NGS: PIK3CA mutation detected

Expression: Luminal A signature

Integrated Radiopathologic-Molecular Diagnosis

## 4 Molecular Profiles

Multi-omic molecular profiling captures the complete molecular state of a patient's disease, spanning from DNA variations to protein expression and metabolic signatures.

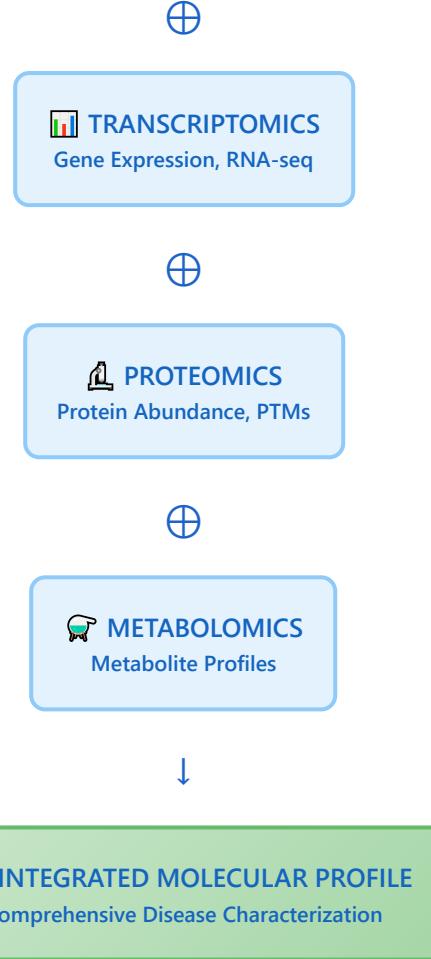
GENOMICS  
Variants, CNVs, SVs

## Omics Data Layers

- **Genomics:** Germline and somatic mutations, CNVs, structural variants (WGS/WES)
- **Transcriptomics:** Gene expression levels, splice variants, fusion transcripts (RNA-seq)
- **Epigenomics:** DNA methylation, histone modifications, chromatin accessibility
- **Proteomics:** Protein abundance, post-translational modifications
- **Metabolomics:** Small molecule metabolites and metabolic pathway activity

## Integration Framework

- Multi-omics network analysis to identify disease drivers
- Pathway enrichment across molecular layers
- Molecular subtype classification for precision therapy
- Biomarker discovery through integrated analysis



## 5 Temporal Alignment

Temporal alignment is critical for understanding disease progression, treatment response, and the dynamic interplay between clinical and molecular changes over time.

Day 0: Diagnosis | Baseline Labs | Initial Genomic

## Temporal Data Challenges

- **Asynchronous Sampling:** Clinical labs, imaging, and molecular assays performed at different times
- **Variable Intervals:** Irregular follow-up schedules and missing data points
- **Treatment Effects:** Therapy-induced changes complicate temporal patterns
- **Biological Lags:** Time delays between molecular changes and clinical manifestations

## Alignment Strategies

- Dynamic time warping for irregular time series
- Interpolation methods for missing timepoints
- Landmark analysis anchored to clinical events (diagnosis, treatment start)
- Longitudinal modeling with mixed-effects approaches
- Causal inference to separate temporal correlation from causation

Profiling

Day 14: Treatment Start | CBC, CMP | Baseline ctDNA

Day 30: Follow-up Labs | CT Scan | ctDNA Monitoring

Day 60: Response Assessment | PET Scan | RNA-seq Analysis

Day 90: Labs | Tumor Markers | Molecular Response Profile

Outcome: Complete Response | Normalized Labs | MRD Negative

# Temporal Integration in Multi-Omics Data Analysis

## Longitudinal Designs

Repeated measurements over time

## Time Series Alignment

Synchronizing different measurement schedules

## Dynamic Modeling

Capturing temporal dynamics

## State Transitions

Disease progression and treatment response

## Trajectory Inference

Reconstructing continuous processes

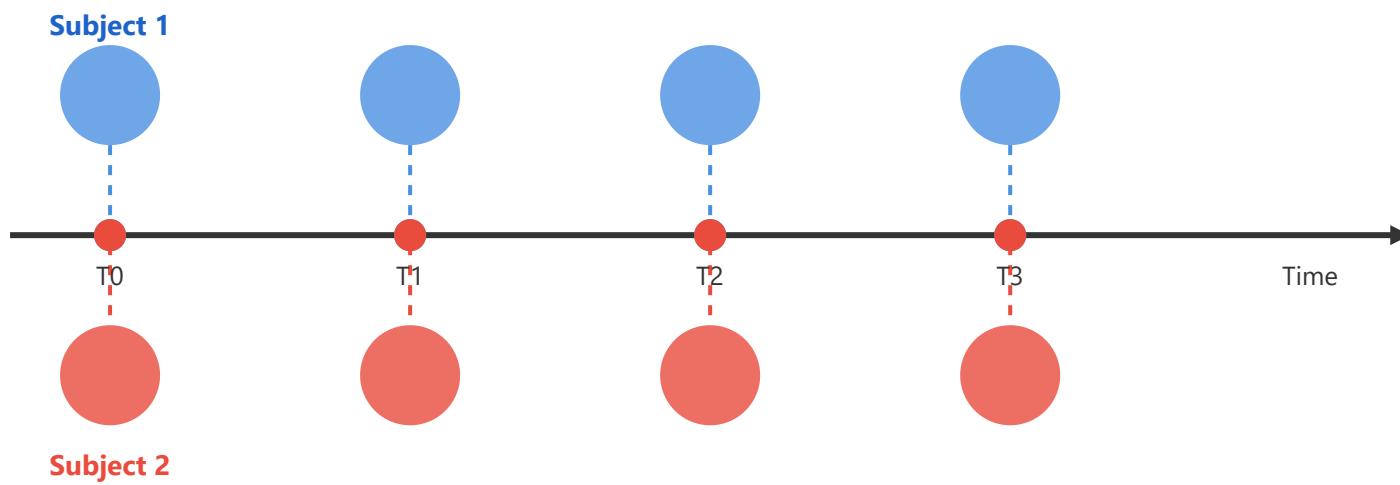
1

## Longitudinal Designs

Longitudinal designs involve collecting multiple samples from the same individuals at different time points. This approach enables researchers to track biological changes within subjects over time, reducing inter-individual variability and increasing statistical power to detect temporal patterns. Unlike cross-sectional studies that provide snapshots at single time points, longitudinal designs capture the dynamic nature of biological processes.

Longitudinal Study Design Visualization

## Repeated Measurements Over Time



### Key Advantages

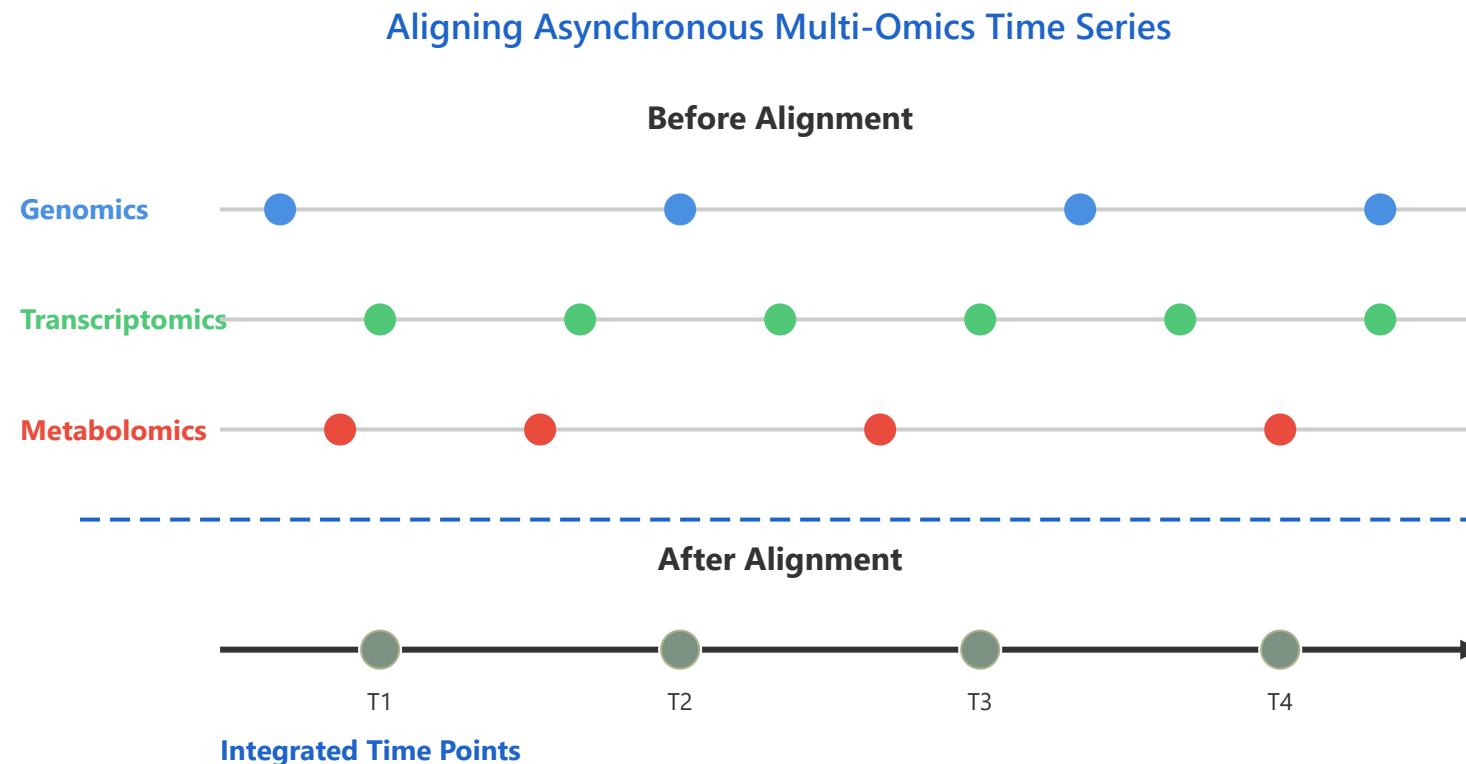
- ▶ Reduced inter-individual variability by using subjects as their own controls
- ▶ Ability to detect individual-specific temporal patterns and trajectories
- ▶ Enhanced statistical power for identifying time-dependent changes
- ▶ Capability to distinguish age effects from cohort effects

### Clinical Example:

Tracking biomarker changes in cancer patients before treatment, during chemotherapy, and post-treatment to understand individual response patterns and predict treatment outcomes.

## 2 Time Series Alignment

Time series alignment addresses the challenge of synchronizing multi-omics data collected at different time intervals or schedules. Different omics layers often have varying measurement frequencies due to technical constraints, cost considerations, or clinical protocols. Alignment methods enable the integration of these asynchronous data streams to create a unified temporal framework for analysis.



### Alignment Strategies

- ▶ Interpolation methods for estimating missing time points
- ▶ Dynamic Time Warping (DTW) to find optimal alignment between sequences
- ▶ Reference time point selection based on clinical events
- ▶ Statistical models accounting for measurement timing uncertainty

#### **Application Example:**

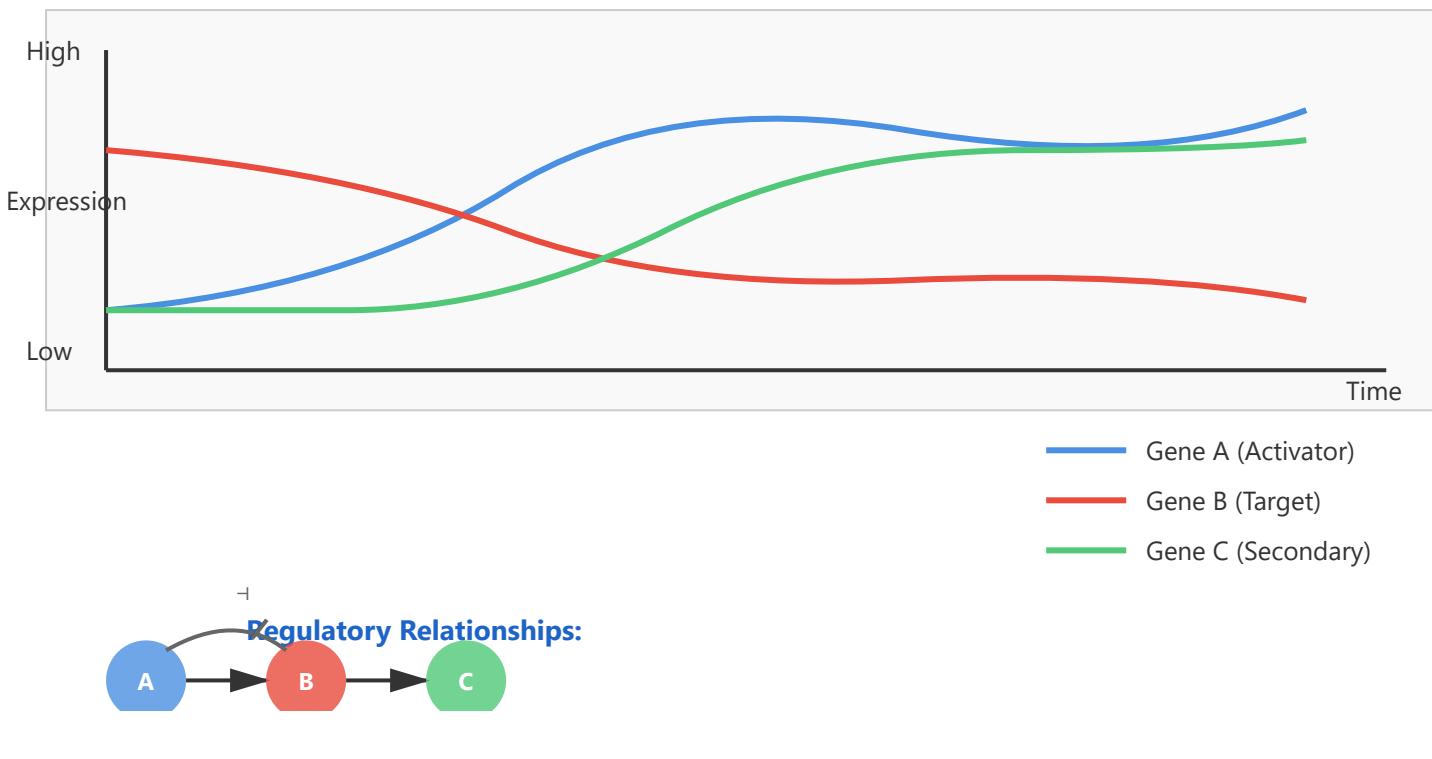
Aligning daily metabolomics measurements with weekly transcriptomics and monthly genomics data in a diabetes study to identify coordinated molecular changes associated with disease progression.

## **3 Dynamic Modeling**

Dynamic modeling captures the temporal evolution of biological systems by representing how molecular states change over time and interact with each other. These models incorporate the rates of change, feedback loops, and regulatory mechanisms that govern biological processes. Common approaches include ordinary differential equations (ODEs), state-space models, and dynamic Bayesian networks.

#### **Dynamic System Modeling Example**

## Gene Regulatory Network Dynamics



### Modeling Approaches

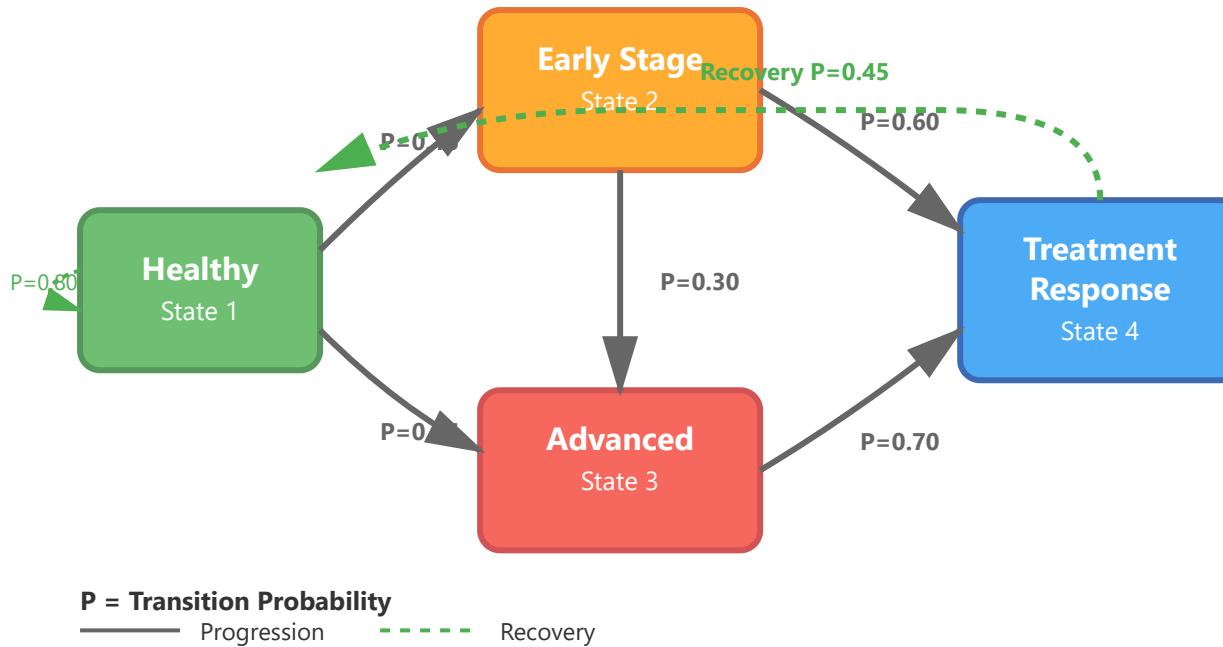
- ▶ Ordinary Differential Equations (ODEs) for continuous dynamics
- ▶ State-Space Models for hidden state estimation
- ▶ Dynamic Bayesian Networks for probabilistic inference
- ▶ Vector Autoregressive (VAR) models for multivariate time series

### Research Application:

## 4 State Transitions

State transition analysis focuses on identifying and characterizing discrete biological states and the transitions between them. This approach is particularly valuable for understanding disease progression, treatment response, and developmental stages. Markov models, Hidden Markov Models (HMMs), and finite state machines are commonly used to represent these transitions and estimate transition probabilities.

### Disease Progression State Transitions



## Key Concepts

- Discrete state identification using clustering or classification methods
- Transition probability estimation from longitudinal observations
- Hidden Markov Models for inferring unobserved states
- Multi-state models for complex transition pathways

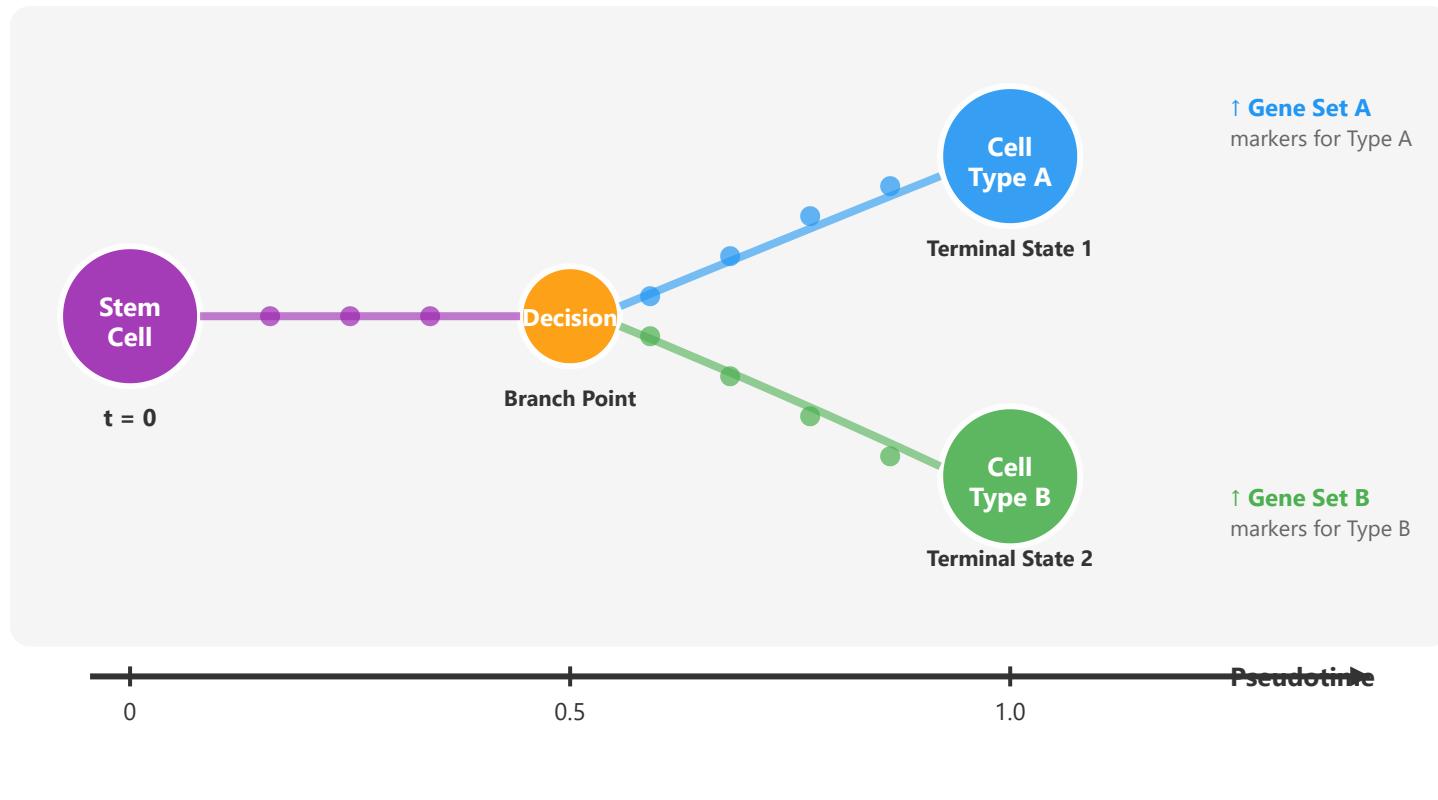
## Clinical Application:

## 5 Trajectory Inference

Trajectory inference reconstructs continuous biological processes from snapshot data, ordering samples along pseudotime to reveal developmental or disease progression pathways. This approach is especially powerful when working with single-cell data or when true temporal information is unavailable. Methods construct trajectories in high-dimensional molecular space, identifying branch points, decision nodes, and terminal states.

### Trajectory Inference and Pseudotime Ordering

## Cell Differentiation Trajectory



### Inference Methods

- ▶ Principal curve and manifold learning algorithms
- ▶ Minimum spanning tree approaches for trajectory construction
- ▶ RNA velocity for directional inference in single-cell data
- ▶ Diffusion pseudotime for complex branching structures

### Single-Cell Application:

Reconstructing hematopoietic differentiation trajectories from single-cell multi-omics data to identify transcription factors and epigenetic modifications that drive cell fate decisions at branch points.

## Spatial Multi-Omics Integration - Detailed Guide

### Five Key Categories of Spatial Integration

#### 1. Spatial Omics

---

Spatial omics technologies enable the measurement of molecular features while preserving their spatial context within tissue samples. This revolutionary approach combines traditional omics profiling with spatial information, allowing researchers to understand not just what molecules are present, but precisely where they are located within the tissue architecture.

## Spatial Omics Technologies

### Spatial Transcriptomics



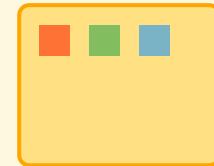
- 10x Visium, MERFISH
- 55µm resolution
- ~20,000 genes

### Spatial Proteomics



- CODEX, IMC
- 1µm resolution
- 40-100 proteins

### Spatial Metabolomics



- MALDI-MSI, DESI-MSI
- 10-50µm resolution
- 100-1000s metabolites

## Key Technologies and Characteristics:

- ▶ **Spatial Transcriptomics:** Technologies like 10x Visium and MERFISH capture gene expression patterns across tissue sections with spatial barcodes or fluorescent probes, enabling whole transcriptome profiling at cellular or subcellular resolution
- ▶ **Spatial Proteomics:** Methods such as CODEX and Imaging Mass Cytometry (IMC) use multiplexed antibody staining to simultaneously detect dozens of proteins while maintaining spatial relationships between cells
- ▶ **Spatial Metabolomics:** Mass spectrometry imaging (MALDI-MSI, DESI-MSI) maps the distribution of metabolites across tissue sections, revealing metabolic heterogeneity and functional zonation
- ▶ **Resolution Trade-offs:** Higher molecular coverage often comes at the cost of spatial resolution, requiring careful technology selection based on research questions
- ▶ **Data Integration:** Combining multiple spatial omics modalities on serial or adjacent tissue sections provides comprehensive molecular characterization

### Clinical and Research Applications:

Spatial omics enables tumor microenvironment mapping to identify therapy resistance mechanisms, neuroscience studies of regional brain heterogeneity, developmental biology investigations of tissue organization during embryogenesis, and precision medicine approaches for patient stratification based on spatial molecular signatures.

## 2. Image Registration

Image registration is the computational process of spatially aligning multiple images or datasets from different sources, modalities, or time points into a common coordinate system. In spatial multi-omics, this enables the integration of complementary information from various profiling techniques while preserving spatial relationships.

### Registration Approaches and Considerations:

- ▶ **Feature-Based Methods:** Detect and match distinctive landmarks (nuclei, vessels, tissue boundaries) between images using algorithms like SIFT, SURF, or deep learning feature extractors
- ▶ **Intensity-Based Methods:** Optimize similarity metrics (mutual information, correlation) directly from image intensities without explicit feature detection
- ▶ **Transformation Models:** Rigid transformations preserve shapes and distances, affine transformations allow scaling and shearing, while non-rigid deformations accommodate tissue distortions from sectioning and processing
- ▶ **Multi-Resolution Strategies:** Pyramid approaches register images at progressively finer scales, improving robustness and computational efficiency

- ▶ **Validation:** Assess registration quality using metrics like target registration error, Dice coefficient for segmentation overlap, and visual inspection of anatomical landmarks

### Practical Challenges and Solutions:

Image registration in spatial multi-omics faces challenges including tissue distortion from sectioning and mounting, differences in image modalities, varying resolutions and fields of view, and computational demands for large datasets. Solutions include using fiducial markers embedded in tissue, applying physics-based deformation models, and leveraging GPU acceleration and cloud computing.

## 3. Cellular Neighborhoods

Cellular neighborhoods represent the local microenvironment surrounding each cell, defined by the identity, spatial arrangement, and molecular states of neighboring cells. Understanding these neighborhoods is crucial because cells do not exist in isolation; their behavior and fate are profoundly influenced by their neighbors through direct cell-cell interactions, paracrine signaling, and shared access to local resources.

### Key Concepts and Methodologies:

- ▶ **Neighborhood Definition:** Defined by radius-based approaches (all cells within distance  $r$ ), k-nearest neighbors, or graph-based connectivity, with optimal parameters depending on tissue density and cell size
- ▶ **Compositional Analysis:** Quantify the relative proportions of different cell types within each neighborhood to identify recurring spatial patterns and microenvironment subtypes

- ▶ **Spatial Statistics:** Use metrics like Ripley's K function, spatial autocorrelation, and permutation tests to assess whether cell-type distributions are random, clustered, or dispersed
- ▶ **Machine Learning Clustering:** Apply unsupervised methods (k-means, hierarchical clustering, graph-based community detection) to discover neighborhood archetypes across samples
- ▶ **Functional Profiling:** Integrate molecular features (gene expression, protein markers) with neighborhood composition to link spatial organization to biological function

### Biological and Clinical Significance:

Cellular neighborhood analysis has revealed that tumors with high immune cell infiltration respond better to immunotherapy compared to immune-excluded or immune-desert phenotypes. In neurodegenerative diseases, specific neuronal-glial neighborhoods correlate with disease progression. During development, stem cell neighborhoods with particular niche compositions determine differentiation outcomes.

## 4. Tissue Architecture

Tissue architecture refers to the large-scale spatial organization and structural patterns that govern tissue function and homeostasis. Beyond individual cell neighborhoods, tissues exhibit hierarchical organization with distinct anatomical regions, functional zonation, and structural boundaries that coordinate physiological processes.

### Architectural Analysis Approaches:

- ▶ **Hierarchical Decomposition:** Tissues exhibit organization at multiple scales from whole organ (centimeters) to functional units (millimeters) to cellular neighborhoods (micrometers), requiring multi-resolution analysis strategies
- ▶ **Boundary Detection:** Identifying interfaces between tissue compartments using image gradients, machine learning segmentation, or molecular marker transitions reveals functional boundaries and barriers
- ▶ **Spatial Pattern Recognition:** Computational approaches detect recurring structural motifs like tubular networks, layered organization, or follicular structures that correlate with tissue function
- ▶ **Graph-Based Modeling:** Representing tissues as networks where nodes are cells or functional units and edges represent spatial proximity or functional connections enables topological analysis
- ▶ **Morphometric Analysis:** Quantifying shape features (circularity, elongation, tortuosity) of tissue structures provides objective measures of architectural disruption in disease

### Disease Applications and Biomarkers:

Architectural disruption is a hallmark of many diseases and provides diagnostic and prognostic information. In cancer, loss of normal tissue architecture indicates invasive potential. Liver fibrosis progression involves architectural remodeling from normal lobular structure to cirrhotic nodules. Architectural features like glandular organization are incorporated into clinical grading systems because they predict outcomes independently of molecular markers.

## 5. 3D Reconstruction

Three-dimensional reconstruction transforms serial 2D tissue sections into comprehensive 3D models that reveal the true spatial organization of tissues and organs. Since biological structures are inherently three-dimensional, 2D sections can be misleading or

incomplete. 3D reconstruction enables volumetric analysis, accurate spatial relationship quantification, and visualization of complex structures that span multiple sections.

## Technical Approaches and Methodologies:

- ▶ **Serial Section Imaging:** Physical sectioning of tissue at regular intervals (5-10 $\mu$ m) followed by imaging each section, requiring careful tracking of section order and orientation
- ▶ **Registration Pipeline:** Sequentially align sections using rigid, affine, or non-rigid transformations to correct for rotation, translation, and tissue deformation during sectioning
- ▶ **Interpolation Methods:** Estimate tissue structure between sections using linear, spline, or shape-based interpolation to create smooth 3D volumes from discrete slices
- ▶ **Surface Rendering:** Extract 3D surfaces from segmented structures using marching cubes or level set methods, enabling visualization of complex anatomical features
- ▶ **Volume Rendering:** Directly visualize 3D intensity data using ray casting or texture mapping, revealing internal structures without explicit segmentation
- ▶ **Multi-Modal Integration:** Combine 3D reconstructions from different imaging modalities in common coordinate space
- ▶ **Computational Requirements:** High-resolution 3D reconstructions can generate terabyte-scale datasets requiring GPU acceleration and efficient data compression

## Scientific and Clinical Impact:

3D reconstruction has transformed understanding of complex biological structures. In neuroscience, connectomics projects trace neural circuits across entire brain regions. Cancer research benefits from accurate tumor volume measurements for treatment response monitoring. Virtual pathology enables digital sectioning of 3D-reconstructed tissues in any orientation. The integration of spatial omics

with 3D reconstruction creates comprehensive molecular atlases that map gene expression, protein distribution, and metabolite patterns throughout entire organs.

**Part 3/3:**

## **Clinical Applications and Future Directions**

- Disease understanding
- Clinical translation
- Future directions

# Disease Subtyping

Advanced Methods for Identifying and Characterizing Disease Heterogeneity

## Molecular Subtypes

Identifying disease subtypes from multi-omics data

## Clinical Correlates

Linking subtypes to clinical outcomes and treatment response

## Consensus Clustering

Robust subtype identification through ensemble methods

## Stability Analysis

Assessing subtype reproducibility and confidence

## Validation Cohorts

Independent validation across multiple datasets

1

## Molecular Subtypes

Molecular subtyping aims to classify diseases based on underlying molecular characteristics rather than clinical symptoms alone. This approach leverages multi-omics data including genomics, transcriptomics, proteomics, and metabolomics to identify distinct disease subtypes with different biological mechanisms.

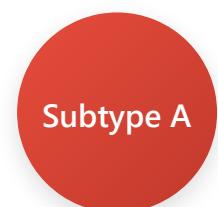
The integration of multiple data types allows for a more comprehensive understanding of disease heterogeneity, leading to more precise diagnosis and personalized treatment strategies.

# Multi-Omics Integration Approach

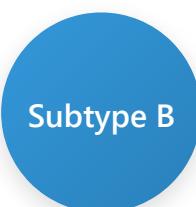
## Multi-Omics Data Integration for Subtype Discovery



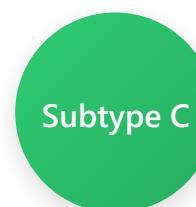
### Example: Cancer Molecular Subtypes



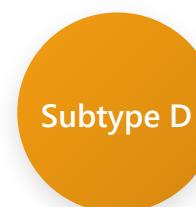
Immune-enriched  
Good prognosis



Proliferative  
Intermediate



Metabolic  
Treatment-resistant



Mesenchymal  
Poor prognosis

### Key Considerations:

- Feature selection is critical - focus on biologically relevant markers
- Batch effects must be corrected across different omics platforms
- Integration methods include early fusion, late fusion, and intermediate fusion

- Dimensionality reduction (PCA, t-SNE, UMAP) helps visualize subtypes

### Real-World Example: Breast Cancer Subtypes

The PAM50 classifier identifies five intrinsic subtypes of breast cancer based on gene expression: Luminal A, Luminal B, HER2-enriched, Basal-like, and Normal-like. Each subtype has distinct molecular characteristics, prognosis, and treatment response patterns.

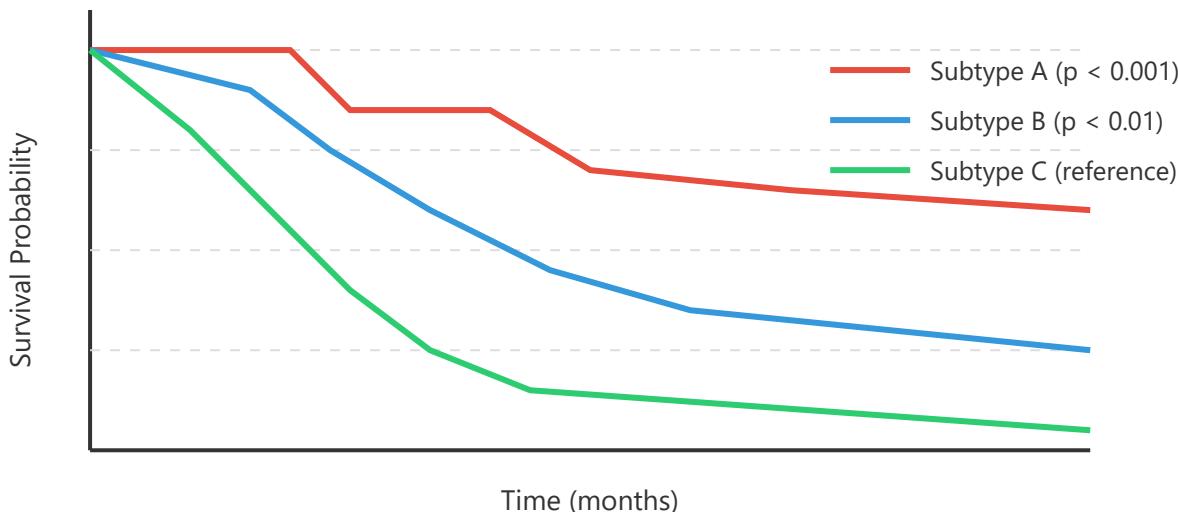
## 2 Clinical Correlates

After identifying molecular subtypes, it's essential to establish their clinical relevance by linking them to patient outcomes, treatment response, and disease progression. This validation ensures that molecular classifications translate into actionable clinical insights.

Clinical correlates include survival outcomes, treatment efficacy, disease recurrence, quality of life measures, and biomarker responses. Statistical methods such as Kaplan-Meier survival analysis, Cox proportional hazards models, and logistic regression are commonly employed.

### Survival Analysis by Subtype

#### Kaplan-Meier Survival Curves by Disease Subtype



## Treatment Response Analysis

**85%**

Subtype A Response Rate

**62%**

Subtype B Response Rate

**38%**

Subtype C Response Rate

**25%**

Subtype D Response Rate

### Statistical Methods:

- Log-rank test for comparing survival curves between subtypes
- Multivariate Cox regression to adjust for confounding variables
- Chi-square tests for categorical outcome associations
- ROC analysis for predictive performance evaluation

### Clinical Translation Example:

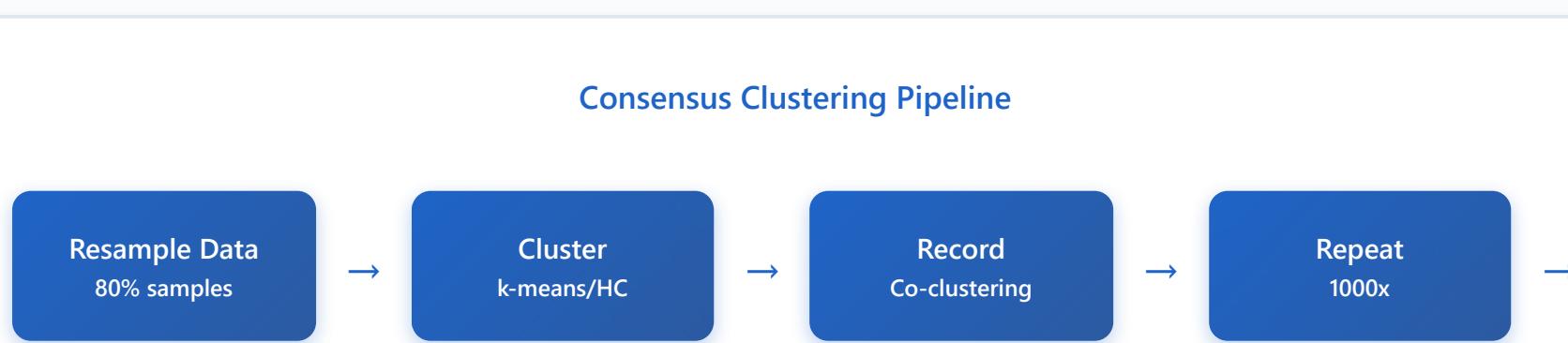
In acute myeloid leukemia (AML), molecular subtypes identified through gene expression profiling showed significantly different responses to standard chemotherapy. Patients with favorable-risk subtypes achieved 70% complete remission rates, while adverse-risk subtypes showed only 30% response, leading to risk-adapted treatment protocols.

### 3 Consensus Clustering

Consensus clustering is a robust method for identifying stable clusters by aggregating results from multiple clustering runs with resampling. This approach addresses the instability inherent in many clustering algorithms and provides a measure of cluster confidence.

The method involves repeatedly subsampling the data, performing clustering on each subsample, and then constructing a consensus matrix that records how frequently pairs of samples cluster together. The final clustering is derived from this consensus matrix.

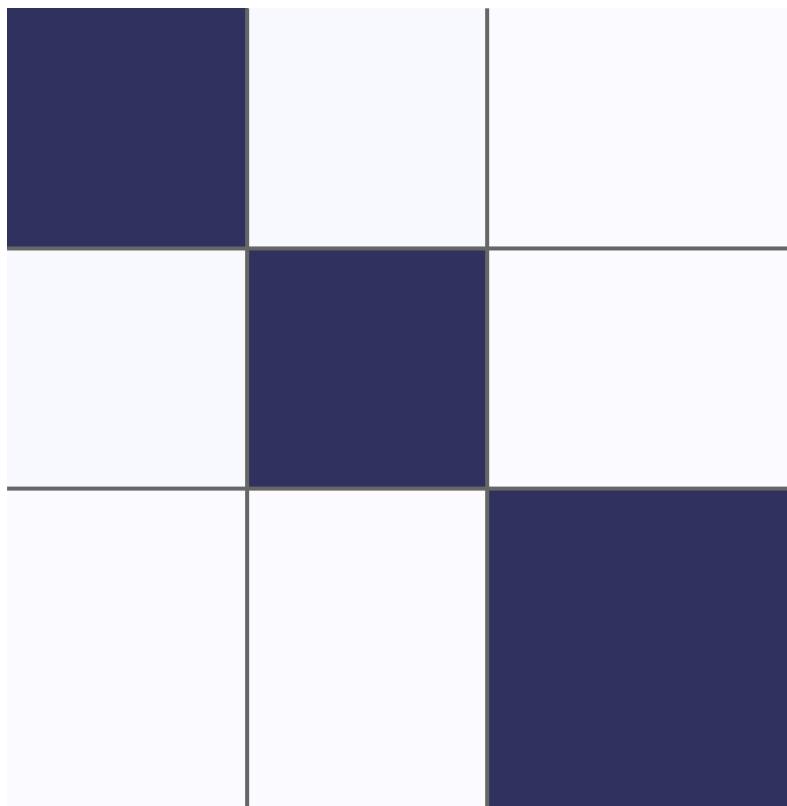
#### Consensus Clustering Workflow



## Consensus Matrix

### Consensus Matrix Heatmap

Example Consensus Matrix (High stability = Dark colors)

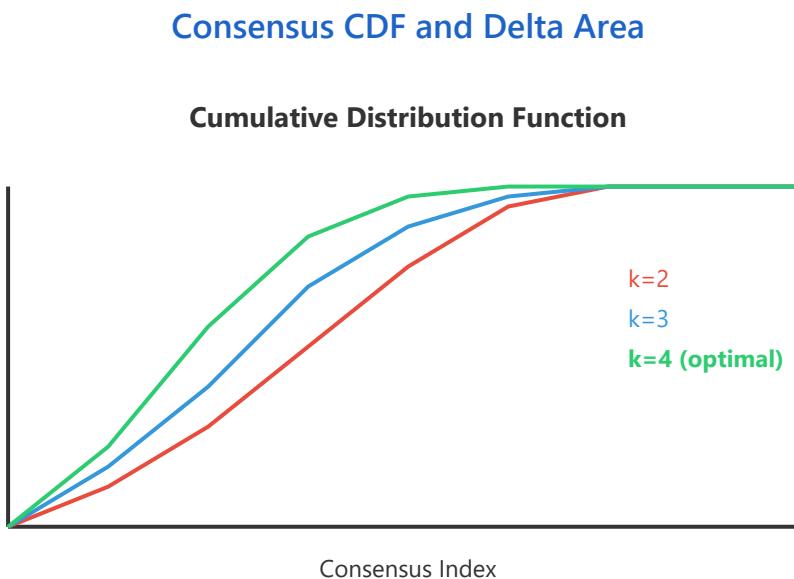


**Interpretation:** Dark diagonal blocks indicate high consensus within clusters. Light off-diagonal regions show low co-clustering between different subtypes.

### Advantages of Consensus Clustering:

- Reduces sensitivity to initialization and outliers
- Provides quantitative measure of cluster stability
- Helps determine optimal number of clusters
- Identifies ambiguous samples with low consensus

## Determining Optimal Cluster Number



### Practical Application:

The Cancer Genome Atlas (TCGA) used consensus clustering to identify robust molecular subtypes across multiple cancer types. For ovarian cancer, they identified four transcriptional subtypes with distinct biological characteristics and clinical outcomes, which have been validated in multiple independent cohorts.

4

## Stability Analysis

Stability analysis assesses the reproducibility and robustness of identified subtypes across different conditions, including data perturbations, feature selections, and algorithmic choices. A stable subtyping solution should be resilient to minor variations in the input data.

Multiple approaches exist for evaluating stability, including bootstrap resampling, cross-validation, subsampling analysis, and noise injection. These methods help distinguish true biological subtypes from artifacts of the clustering algorithm.

### Stability Assessment Methods

#### Stability Evaluation Framework

Bootstrap  
Resampling

Feature  
Subsampling

Noise  
Injection



## Stability Metrics

Jaccard, Rand Index, ARI

### Stability Metrics Visualization

**0.92**

Adjusted Rand Index  
(High Stability)

**0.88**

Jaccard Coefficient  
(High Stability)

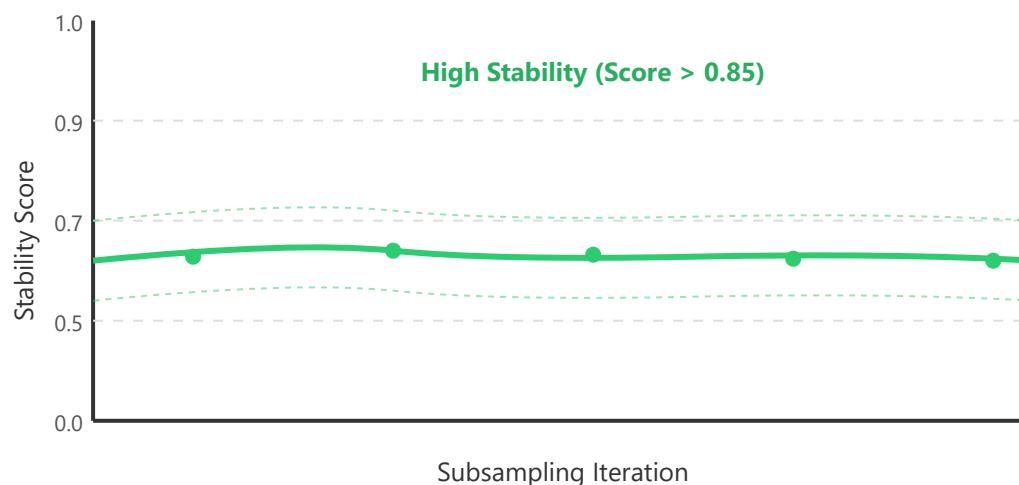
**0.65**

Silhouette Score  
(Moderate)

**3.2**

Gap Statistic  
(k=4 optimal)

### Stability Across Subsampling Iterations



### Stability Metrics Interpretation:

- Adjusted Rand Index (ARI): 0-1 scale, >0.8 indicates high stability
- Jaccard Coefficient: Measures overlap between clustering solutions
- Silhouette Score: Assesses separation between clusters (-1 to 1)
- Stability > 0.85 suggests robust, reproducible subtypes

### Case Study: Alzheimer's Disease Subtypes

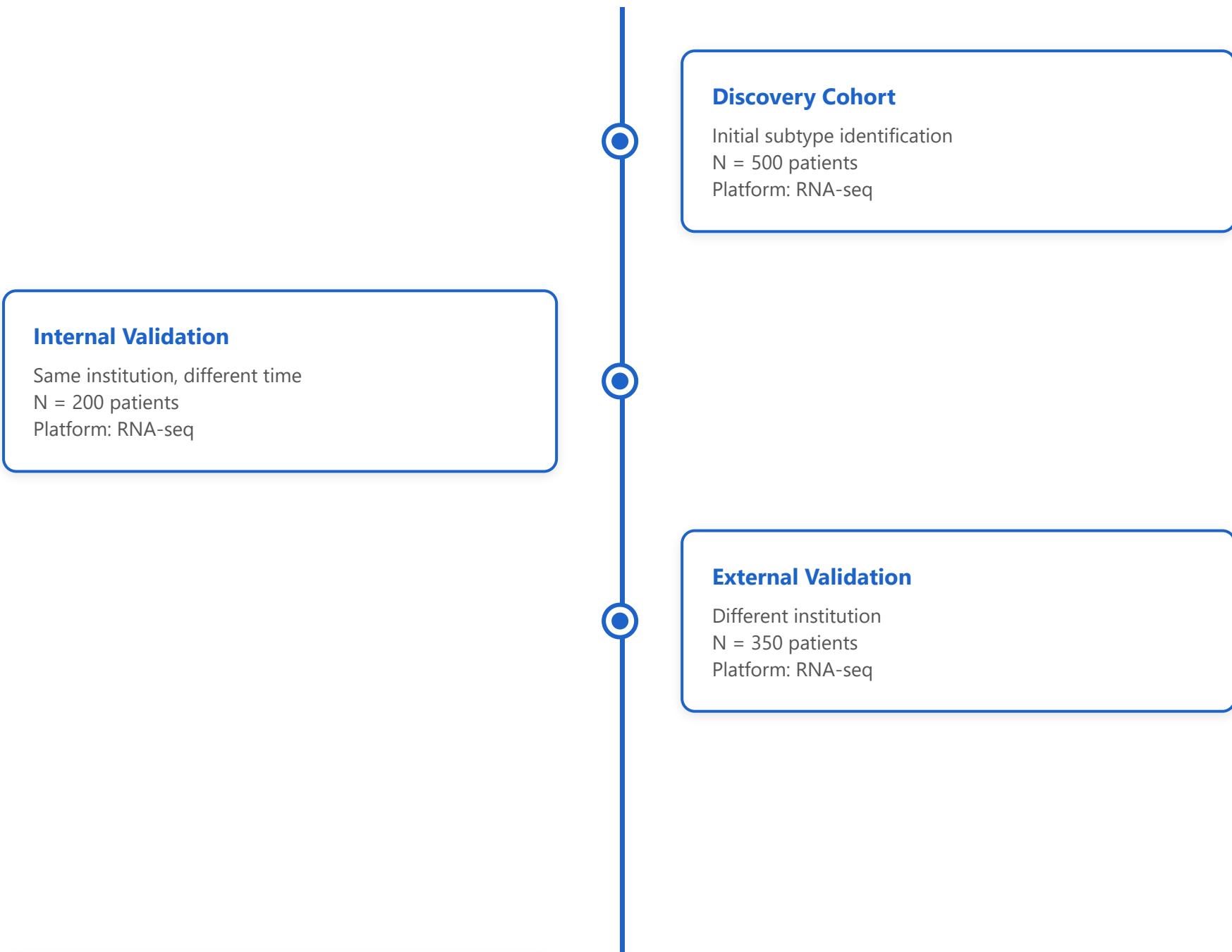
Recent research identified three stable subtypes of Alzheimer's disease through multi-modal neuroimaging and biomarker data. Bootstrap stability analysis with 10,000 iterations showed ARI > 0.90, confirming these subtypes were highly reproducible and not artifacts of the clustering method. Each subtype showed distinct patterns of brain atrophy and cognitive decline trajectories.

## 5 Validation Cohorts

Independent validation is the gold standard for confirming that identified subtypes are generalizable and not specific to the discovery dataset. Validation cohorts should be collected from different populations, time periods, or institutions to ensure broad applicability.

The validation process involves training a classifier on the discovery cohort, then applying it to independent cohorts to assess whether the subtypes maintain their distinct molecular and clinical characteristics. Cross-platform validation is particularly important when different technologies are used.

### Multi-Cohort Validation Framework

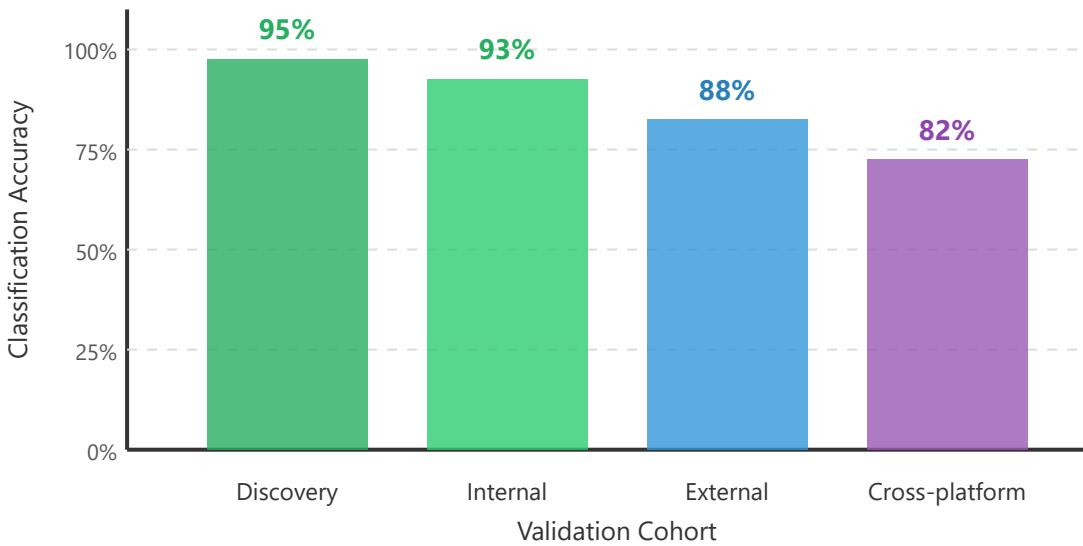


## Cross-Platform Validation

Different technology  
N = 180 patients  
Platform: Microarray

## Validation Performance Metrics

Classification Accuracy Across Validation Cohorts

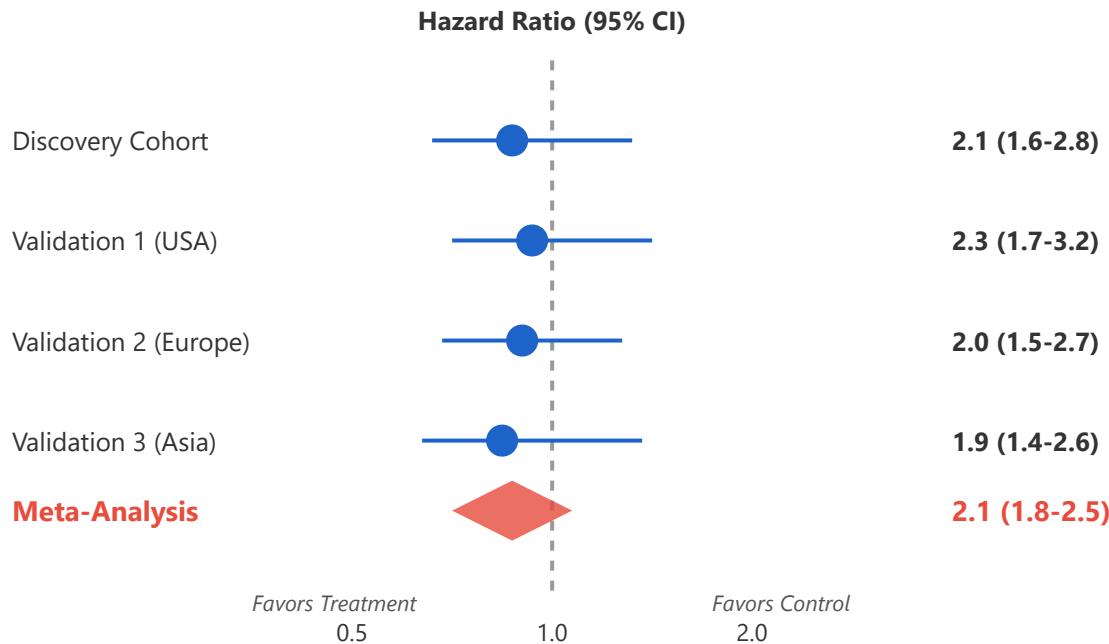


Validation Best Practices:

- Use completely independent datasets not involved in subtype discovery
- Validate across different populations and ethnicities when possible
- Test robustness across different technical platforms
- Confirm both molecular characteristics and clinical associations
- Report confidence intervals and uncertainty estimates

## Multi-Study Meta-Validation

Forest Plot: Hazard Ratios Across Validation Studies



Real-World Validation Success:

The intrinsic subtypes of breast cancer (Basal, HER2-enriched, Luminal A, Luminal B) were initially discovered in a cohort of 85 tumors using hierarchical clustering of gene expression. These subtypes have since been validated in over 50 independent cohorts comprising more than 10,000 patients worldwide, across multiple platforms including microarrays, RNA-seq, and targeted gene panels. The consistency of prognostic associations across all these studies demonstrates the biological and clinical validity of these subtypes.

### Challenges in Validation:

- Batch effects between discovery and validation cohorts must be addressed
- Different platforms may require careful normalization and feature mapping
- Sample size in validation cohorts should be adequate for statistical power
- Publication bias may lead to overestimation of validation success rates



## Summary and Best Practices

Successful disease subtyping requires a systematic approach that integrates molecular data, ensures robustness through consensus methods, validates stability, and confirms findings in independent cohorts. The ultimate goal is to identify clinically actionable subtypes that improve patient stratification and treatment selection.

### Key Takeaways:

1. **Multi-omics integration** provides comprehensive molecular characterization
2. **Clinical correlation** ensures subtypes have practical relevance

3. **Consensus clustering** improves robustness and reliability
4. **Stability analysis** distinguishes real subtypes from artifacts
5. **Independent validation** confirms generalizability across populations

## Future Directions:

Emerging approaches include single-cell multi-omics for higher resolution subtyping, machine learning methods for integrating diverse data types, spatial transcriptomics for tissue architecture analysis, and real-time subtype classification for clinical decision support systems.

# Prognosis Prediction in Precision Medicine

## Multi-modal Signatures

Prognostic signatures from integrated data

## Risk Stratification

Identifying high-risk patients

## Survival Models

Cox regression and deep survival models

## Time-dependent ROC

Evaluating time-to-event predictions

## Clinical Utility

Decision curve analysis

## 1. Multi-modal Signatures

Multi-modal signatures integrate diverse data types to create comprehensive prognostic models that capture the complexity of disease progression. By combining genomic, transcriptomic, proteomic, imaging, and clinical data, these signatures provide a more holistic view of patient prognosis than any single data modality alone.

### Key Components:

- **Genomic Data:** DNA mutations, copy number variations, and structural variants that influence disease outcome
- **Transcriptomic Data:** Gene expression patterns that reflect biological state and treatment response

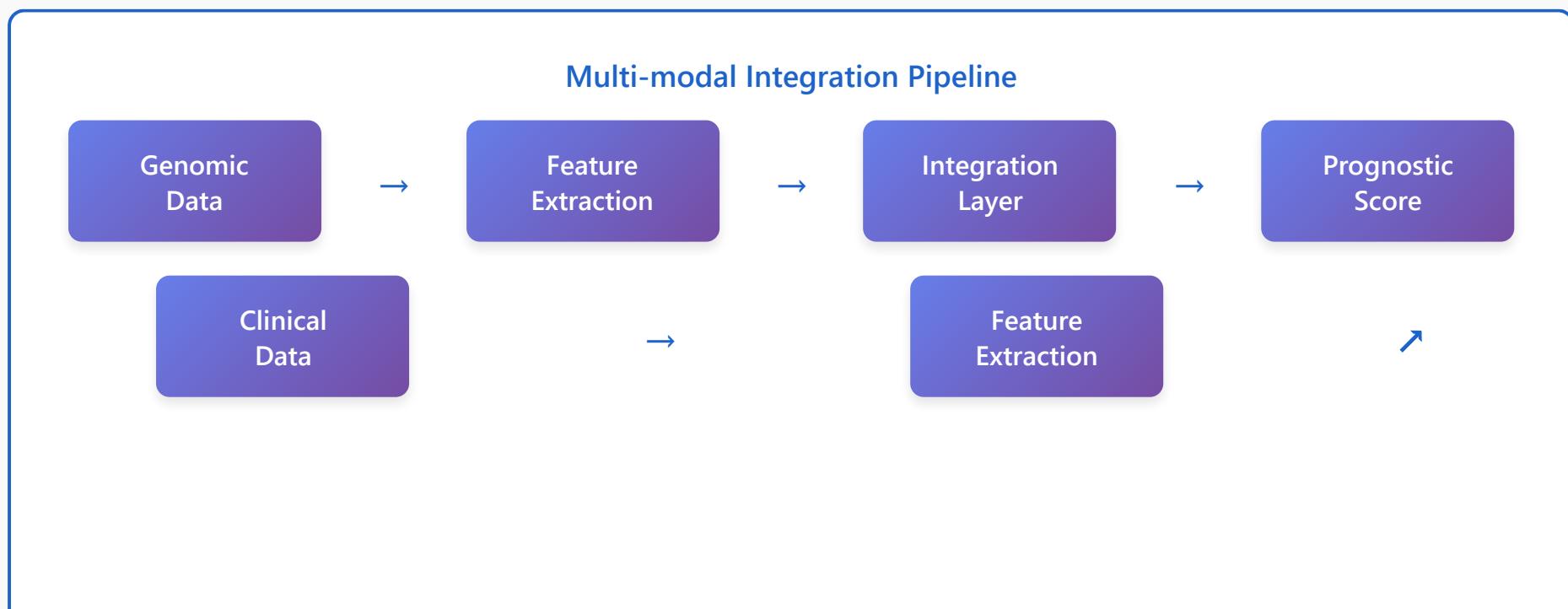
- **Proteomic Data:** Protein abundance and post-translational modifications
- **Clinical Data:** Patient demographics, medical history, and treatment information
- **Imaging Data:** Radiological features and pathology image characteristics

### Integration Approaches:

- **Early Integration:** Concatenate features from all modalities before model training
- **Late Integration:** Train separate models for each modality and combine predictions
- **Intermediate Integration:** Learn shared representations across modalities

#### Clinical Impact:

Multi-modal signatures have shown superior performance in cancer prognosis, often achieving 10-20% improvement in prediction accuracy compared to single-modality approaches.



## 2. Risk Stratification

Risk stratification categorizes patients into distinct groups based on their predicted clinical outcomes, enabling personalized treatment decisions. This approach allows clinicians to identify high-risk patients who may benefit from aggressive interventions and low-risk patients who may safely avoid unnecessary treatments.

### Stratification Methods:

- **Score-based Stratification:** Divide patients using continuous risk scores (e.g., quartiles, tertiles)
- **Clustering Approaches:** Unsupervised grouping based on molecular or clinical profiles
- **Tree-based Methods:** Recursive partitioning to identify natural risk groups
- **Machine Learning Classification:** Supervised learning to predict risk categories

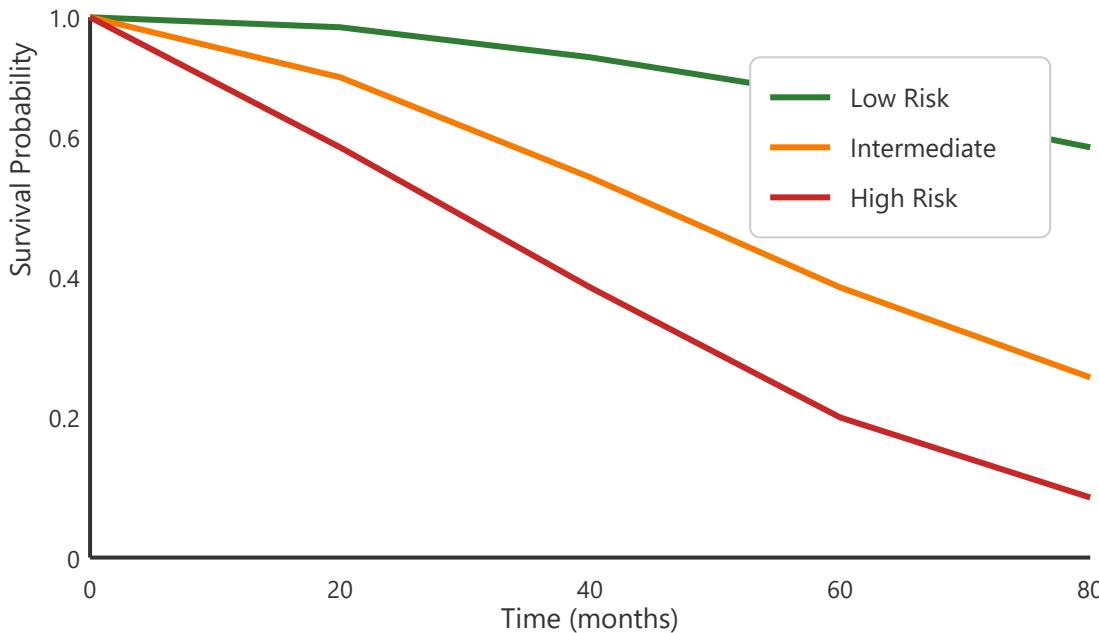
### Common Risk Groups:

- **Low Risk:** Favorable prognosis, may benefit from treatment de-escalation
- **Intermediate Risk:** Standard treatment protocols appropriate
- **High Risk:** Poor prognosis, candidates for intensive therapy or clinical trials

### Example Application:

In breast cancer, the Oncotype DX recurrence score stratifies patients into low, intermediate, and high-risk groups, guiding decisions about adjuvant chemotherapy. Studies show that 70% of patients classified as low-risk can safely avoid chemotherapy.

### Risk Stratification Example



Kaplan-Meier survival curves showing distinct outcomes across risk groups

### 3. Survival Models

Survival models analyze time-to-event data, accounting for censored observations where the event of interest has not occurred by the end of follow-up. These models are essential for prognosis prediction as they handle the temporal nature of clinical outcomes.

#### Cox Proportional Hazards Model:

- **Semi-parametric approach:** Makes no assumptions about baseline hazard function
- **Hazard Ratio:**  $h(t) = h_0(t) \times \exp(\beta_1x_1 + \beta_2x_2 + \dots + \beta_px_p)$
- **Interpretability:** Coefficients represent log-hazard ratios, easily interpretable
- **Limitations:** Assumes proportional hazards, linear relationships

#### Deep Survival Models:

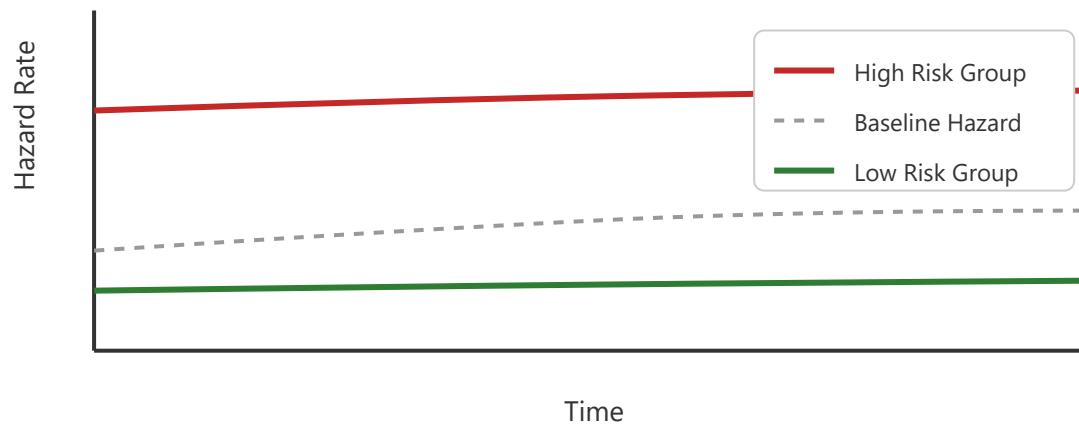
- **DeepSurv:** Neural network extension of Cox model, captures non-linear relationships
- **Neural Multi-Task Logistic Regression:** Predicts discrete-time survival probabilities
- **Variational Autoencoders:** Learn latent representations for survival prediction
- **Advantages:** Handle complex interactions, integrate multiple data types, no proportional hazards assumption

#### Model Selection Guidelines:

Cox regression is preferred when interpretability is crucial and sample sizes are moderate. Deep learning approaches excel with large datasets ( $>1000$  samples) and complex, multi-modal data where non-linear relationships are expected.

Model	Strengths	Use Cases
Cox Regression	Interpretable Statistical rigor	Clinical trials Small datasets
DeepSurv	Non-linear Multi-modal	Large cohorts Complex data
Random Survival Forest	Robust Variable selection	High- dimensional Genomics

### Hazard Function Visualization



## 4. Time-dependent ROC

Time-dependent Receiver Operating Characteristic (ROC) analysis evaluates the discriminative ability of survival models at specific time points. Unlike standard ROC analysis, it accounts for censoring and the time-varying nature of survival predictions.

### Key Concepts:

- **Sensitivity (True Positive Rate):** Proportion of patients who experienced the event by time t and were correctly classified as high-risk
- **Specificity (True Negative Rate):** Proportion of patients who remained event-free by time t and were correctly classified as low-risk
- **AUC(t):** Area under the time-dependent ROC curve at time t; ranges from 0.5 (no discrimination) to 1.0 (perfect discrimination)
- **Time-varying Performance:** Model performance may change over time, requiring evaluation at multiple clinically relevant time points

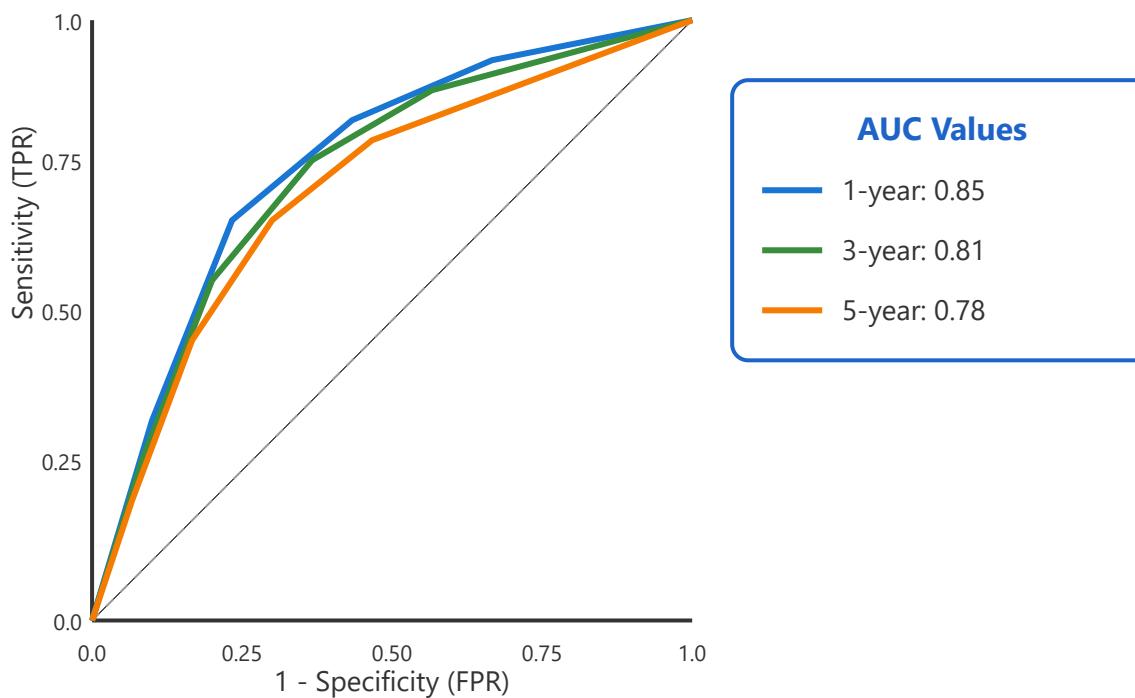
### Evaluation Strategy:

- **Multiple Time Points:** Assess performance at 1-year, 3-year, 5-year survival
- **Integrated AUC:** Summary measure across time range
- **Censoring Adjustments:** Use inverse probability weighting or other methods to handle censored data

### Performance Benchmarks:

In oncology, an  $AUC > 0.70$  is considered acceptable,  $> 0.75$  good, and  $> 0.80$  excellent for prognostic models. Clinical deployment typically requires  $AUC > 0.75$  with validation in independent cohorts.

## Time-dependent ROC Analysis



**Interpretation:** This example shows decreasing discriminative ability over time, which is common in prognostic models. The model performs best for short-term predictions (1-year AUC = 0.85) and slightly worse for long-term predictions (5-year AUC = 0.78), though all remain in the "good" range.

## 5. Clinical Utility

Clinical utility assessment determines whether a prognostic model provides actionable information that improves patient outcomes or clinical decision-making. Decision curve analysis (DCA) is the gold standard for evaluating clinical utility by comparing the net benefit of using a model versus default strategies.

### Decision Curve Analysis (DCA):

- **Net Benefit:** Weighs true positives against false positives, accounting for the relative harm of each
- **Threshold Probability:** The probability at which a patient/clinician would opt for treatment
- **Comparison Strategies:** Model vs. "treat all" vs. "treat none"
- **Clinical Interpretation:** Shows the range of threshold probabilities where the model adds value

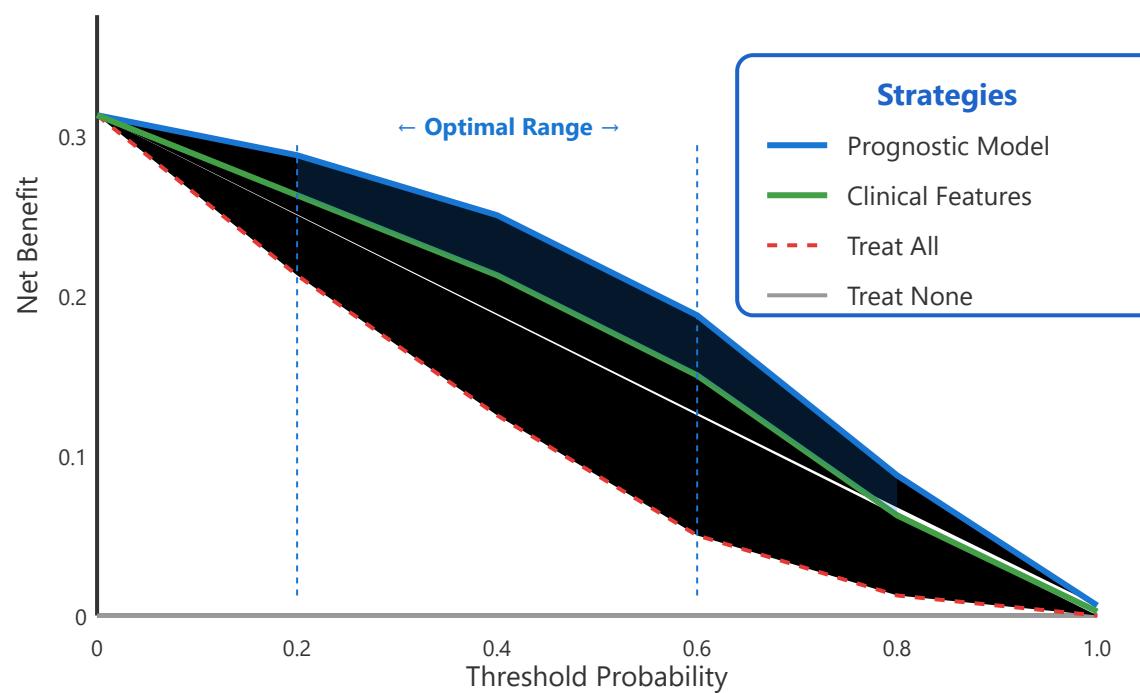
### Components of Clinical Utility:

- **Actionability:** Does the prediction lead to different treatment decisions?
- **Impact on Outcomes:** Does using the model improve patient survival, quality of life, or reduce costs?
- **Implementation Feasibility:** Can the model be integrated into clinical workflow?
- **Cost-effectiveness:** Does the benefit justify the cost of implementation?

#### Real-World Example:

The 21-gene recurrence score in breast cancer demonstrates high clinical utility. Randomized trials showed that using the score to guide treatment decisions resulted in chemotherapy de-escalation in 70% of intermediate-risk patients without compromising survival, while reducing treatment toxicity and healthcare costs by \$2 billion annually in the US.

## Decision Curve Analysis



## Implementation Framework



**Key Insight:** The prognostic model (blue curve) shows superior net benefit compared to treating all patients or using clinical features alone across threshold probabilities of 0.1 to 0.6. This indicates the model is clinically useful for decision-making within this range, where most clinical decisions occur.

## Integration in Clinical Practice

Successful prognostic models combine multi-modal signatures, robust risk stratification, validated survival models, rigorous time-dependent evaluation, and demonstrated clinical utility. The ultimate goal is to provide actionable predictions that improve patient outcomes while being feasible to implement in real-world clinical settings.

## Drug Response Prediction

### Sensitivity Prediction

Predicting drug effectiveness

### Resistance Markers

Identifying resistance mechanisms

### Combination Effects

Drug synergy and antagonism

### Pharmacogenomics

Genetic variants affecting drug response

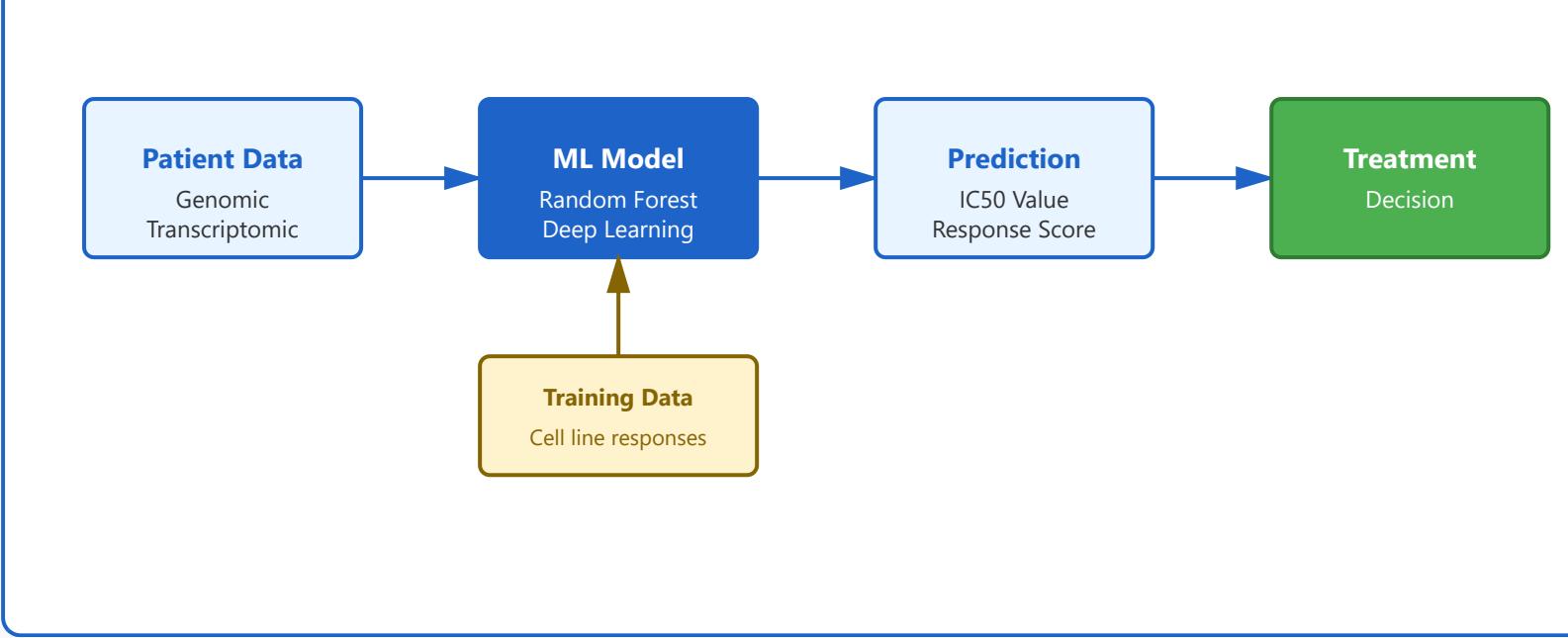
### Clinical Trials

Integration in precision medicine trials

## 1. Sensitivity Prediction

Drug sensitivity prediction utilizes machine learning algorithms to forecast how effectively a particular drug will work against cancer cells or pathogens based on genomic, transcriptomic, and proteomic data. This approach enables personalized treatment selection by identifying patients most likely to benefit from specific therapeutics.

### Drug Sensitivity Prediction Workflow

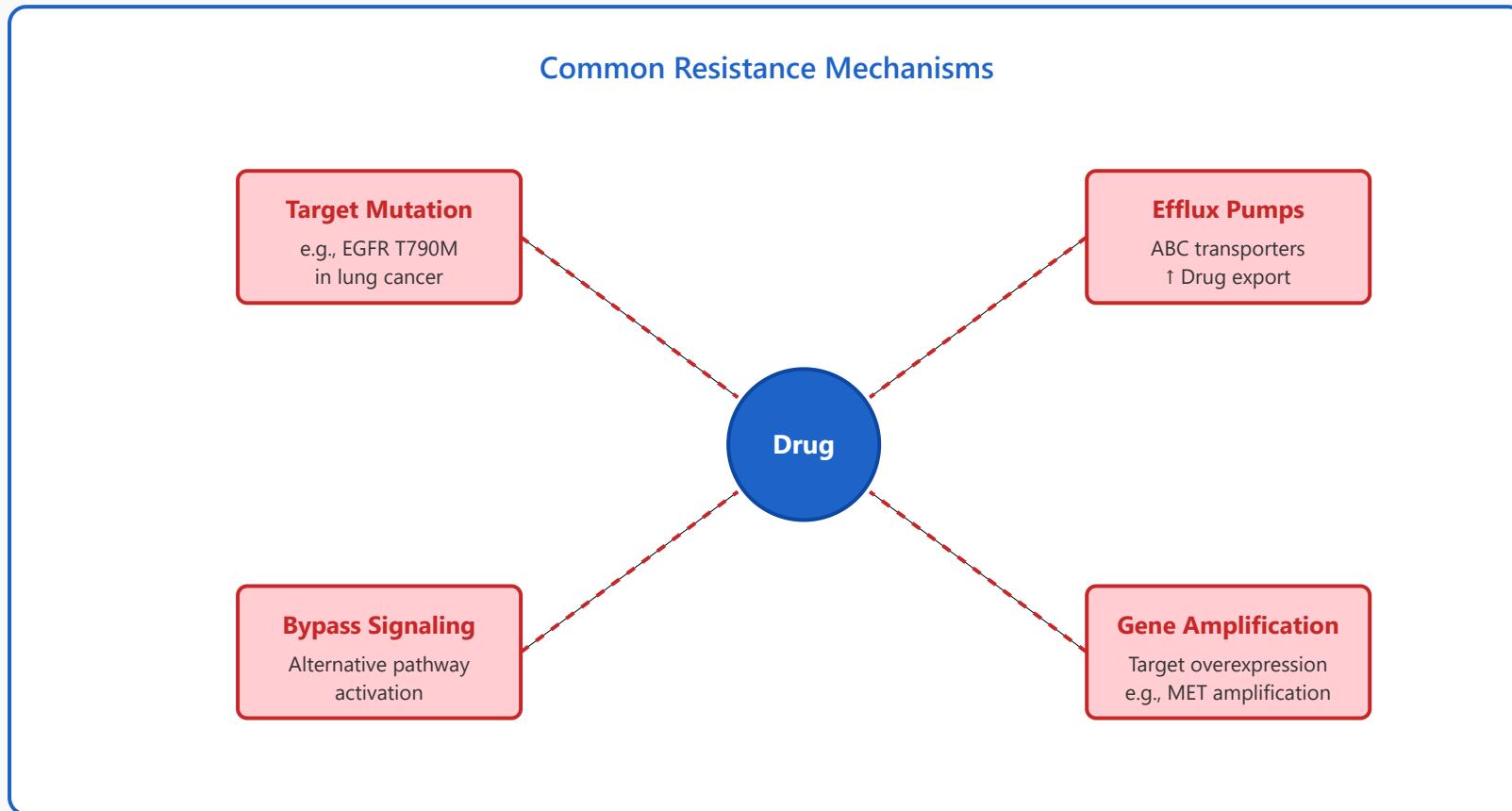


### Key Applications:

- **Cancer Treatment Selection:** Predicting response to chemotherapy, targeted therapy, and immunotherapy based on tumor molecular profiles
- **IC50 Prediction:** Estimating the half-maximal inhibitory concentration to determine optimal drug dosing
- **Biomarker Discovery:** Identifying genomic features (mutations, gene expression patterns) that correlate with drug sensitivity
- **Patient Stratification:** Classifying patients into responder and non-responder groups before treatment initiation

## 2. Resistance Markers

Resistance markers are molecular indicators that signal a patient's tumor or pathogen may not respond to specific drugs. Understanding resistance mechanisms enables clinicians to avoid ineffective treatments and select alternative therapeutic strategies. These markers can be genetic mutations, gene amplifications, or expression changes that confer drug resistance.



### Clinical Examples:

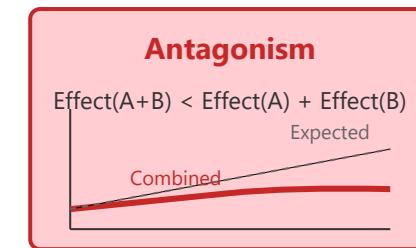
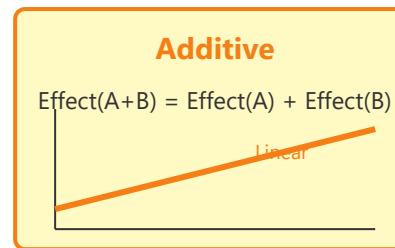
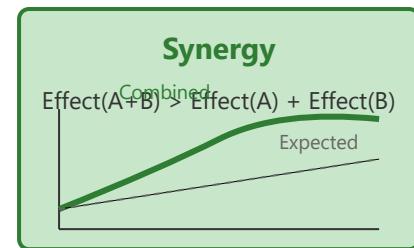
- **EGFR T790M Mutation:** Confers resistance to first-generation EGFR inhibitors in non-small cell lung cancer; detectable through liquid biopsy
- **BCR-ABL Mutations:** Multiple point mutations in chronic myeloid leukemia lead to resistance against tyrosine kinase inhibitors like imatinib

- **MDR1/P-glycoprotein Overexpression:** Increased efflux pump activity reduces intracellular drug concentration in various cancers
- **MSI Status:** Microsatellite instability can predict resistance to certain chemotherapies but sensitivity to immunotherapy
- **Adaptive Resistance:** Dynamic monitoring of emerging resistance markers during treatment enables early intervention strategies

### 3. Combination Effects

Drug combination therapy aims to enhance treatment efficacy through synergistic interactions while minimizing toxicity. Computational prediction of drug synergy and antagonism helps identify optimal combination regimens, reduce trial-and-error in clinical practice, and accelerate drug development. Understanding interaction mechanisms is crucial for rational combination therapy design.

#### Drug Interaction Types



**Drug Combination Matrix**

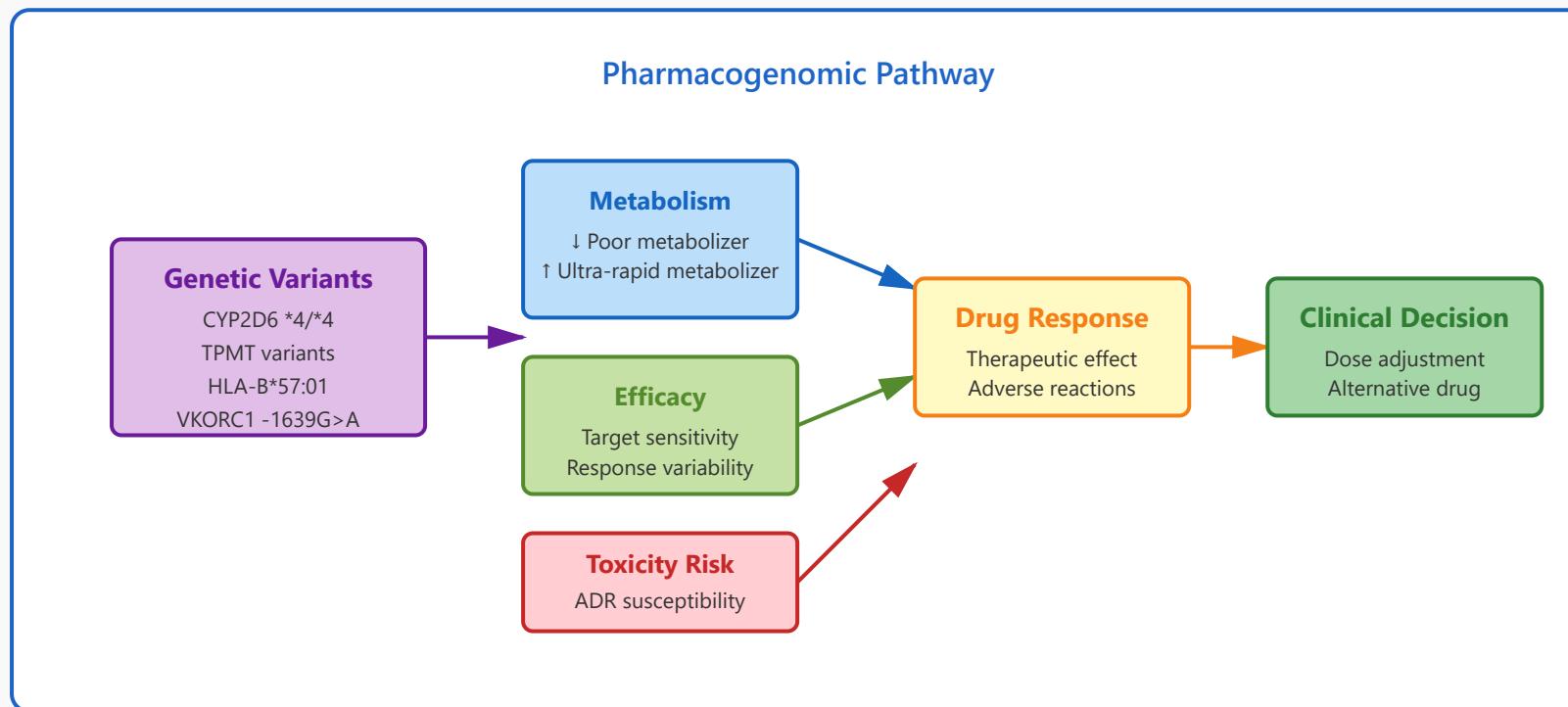
Drug A	+++	++	+
(dose)	++	+++	++

### Combination Strategies:

- **Mechanistic Synergy:** Drugs targeting complementary pathways (e.g., BRAF + MEK inhibitors in melanoma) produce enhanced effects
- **Sequential Blockade:** Combining drugs that prevent compensatory resistance mechanisms (e.g., dual HER2 blockade with trastuzumab and pertuzumab)
- **Bliss Independence Model:** Mathematical framework for quantifying synergy by comparing observed vs. expected effects
- **Network-Based Prediction:** Using systems biology approaches to identify optimal drug combinations based on pathway crosstalk
- **High-Throughput Screening:** Large-scale testing of drug combinations combined with ML to predict synergistic pairs

## 4. Pharmacogenomics

Pharmacogenomics studies how genetic variations affect individual responses to medications, encompassing drug metabolism, efficacy, and adverse reactions. By integrating genomic data with clinical pharmacology, this field enables truly personalized medicine where drug selection and dosing are optimized based on a patient's genetic profile, improving outcomes while reducing toxicity.



### Clinical Applications:

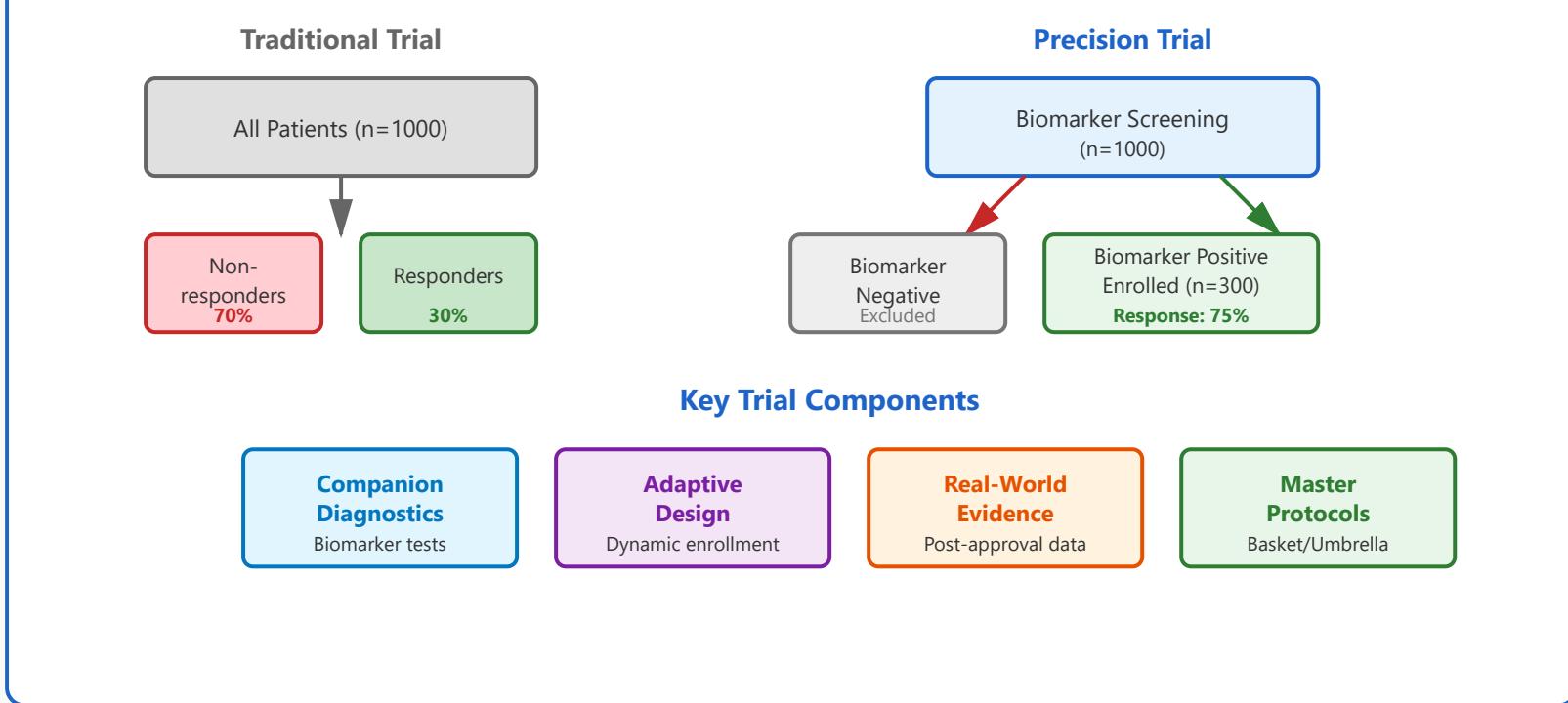
- **CYP450 Polymorphisms:** CYP2D6, CYP2C19, and CYP2C9 variants dramatically affect metabolism of ~25% of commonly prescribed drugs, guiding dosing for antidepressants, anticoagulants, and pain medications
- **Warfarin Dosing:** VKORC1 and CYP2C9 genotypes enable precise warfarin dosing algorithms, reducing bleeding risks and improving anticoagulation control

- **HLA-Based ADR Prevention:** HLA-B\*57:01 screening prevents abacavir hypersensitivity; HLA-B\*15:02 testing avoids carbamazepine-induced Stevens-Johnson syndrome
- **Thiopurine Methyltransferase (TPMT):** Testing identifies patients at high risk for severe myelosuppression from azathioprine or mercaptopurine
- **Oncology Applications:** DPYD variants predict fluoropyrimidine toxicity; UGT1A1\*28 guides irinotecan dosing

## 5. Clinical Trials

Integration of drug response prediction into clinical trial design represents a paradigm shift toward precision medicine. By leveraging biomarkers, genomic data, and predictive algorithms, modern trials can identify patient populations most likely to benefit from experimental therapies, employ adaptive designs, and accelerate drug approval processes while improving success rates and reducing costs.

### Precision Medicine Clinical Trial Design



### Modern Trial Innovations:

- **Basket Trials**: Enroll patients with the same biomarker across different tumor types (e.g., NCI-MATCH, targeting specific mutations regardless of cancer origin)
- **Umbrella Trials**: Test multiple targeted therapies in a single disease type based on different biomarkers (e.g., Lung-MAP for non-small cell lung cancer)
- **Adaptive Trial Designs**: Allow protocol modifications based on interim results, improving efficiency and ethical patient allocation
- **Companion Diagnostics**: Co-develop diagnostic tests with therapeutics to identify appropriate patient populations (e.g., PD-L1 testing for checkpoint inhibitors)
- **Digital Biomarkers**: Incorporate wearable devices and remote monitoring to track real-time response and adverse events

- **Platform Trials:** Perpetual trials that can add or drop treatment arms based on performance, accelerating multiple drug evaluations simultaneously

## Biomarker Panels

### Multi-analyte Tests

Combining multiple biomarkers

### Optimal Combinations

Feature selection for panels

### Performance Metrics

Sensitivity, specificity, PPV, NPV

### Cost-benefit

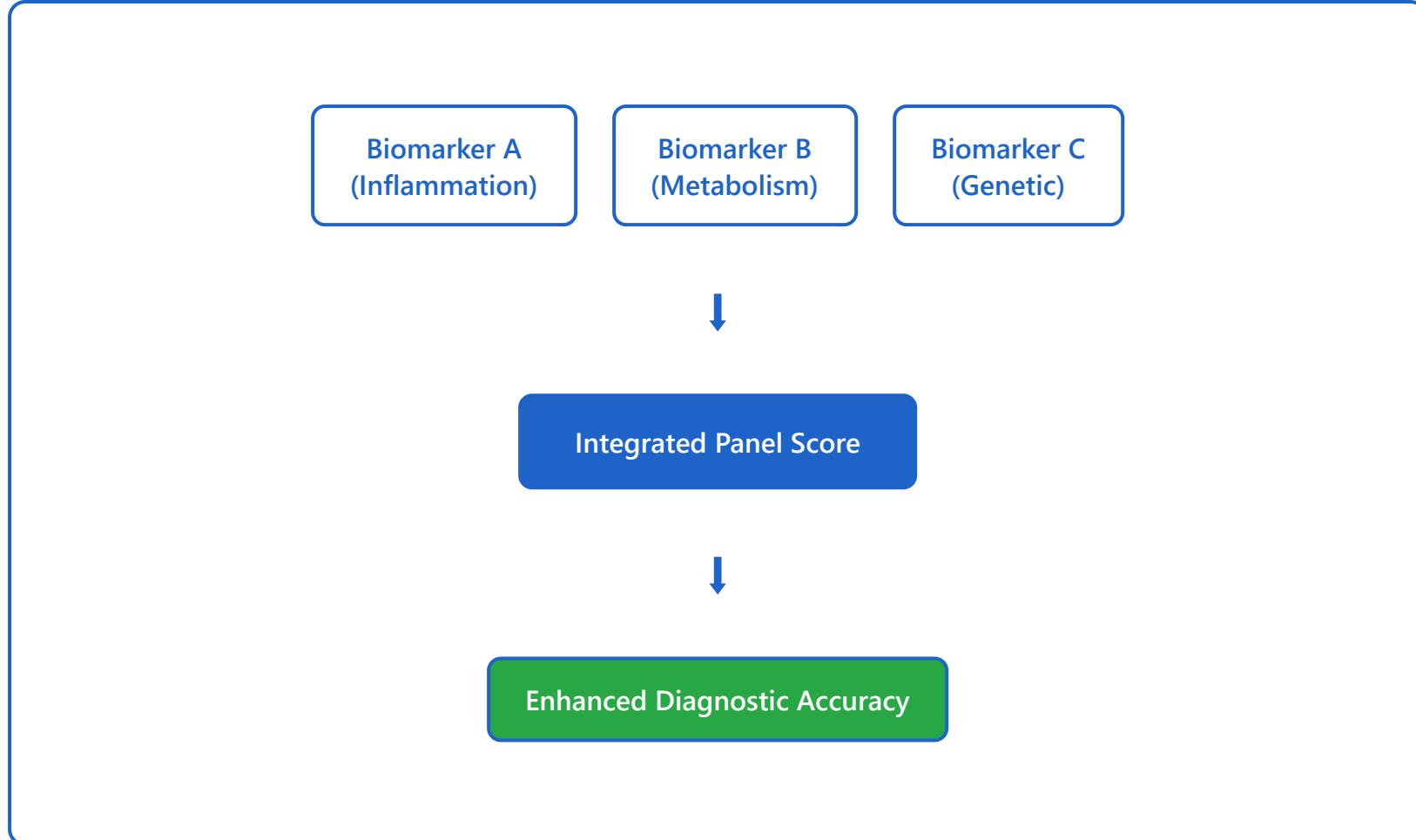
Clinical and economic considerations

### Regulatory Approval

FDA/EMA approval pathways

## 1. Multi-analyte Tests

Multi-analyte tests combine multiple biomarkers into a single assay to improve diagnostic accuracy and provide comprehensive disease assessment. This approach leverages complementary information from different biological pathways to enhance clinical decision-making.



#### Clinical Example: Cardiac Risk Panel

**Traditional:** Single marker (Total Cholesterol)

**Multi-analyte Panel:** LDL, HDL, Triglycerides, hsCRP, Troponin, BNP

**Benefit:** Improved risk stratification from 60% to 85% accuracy

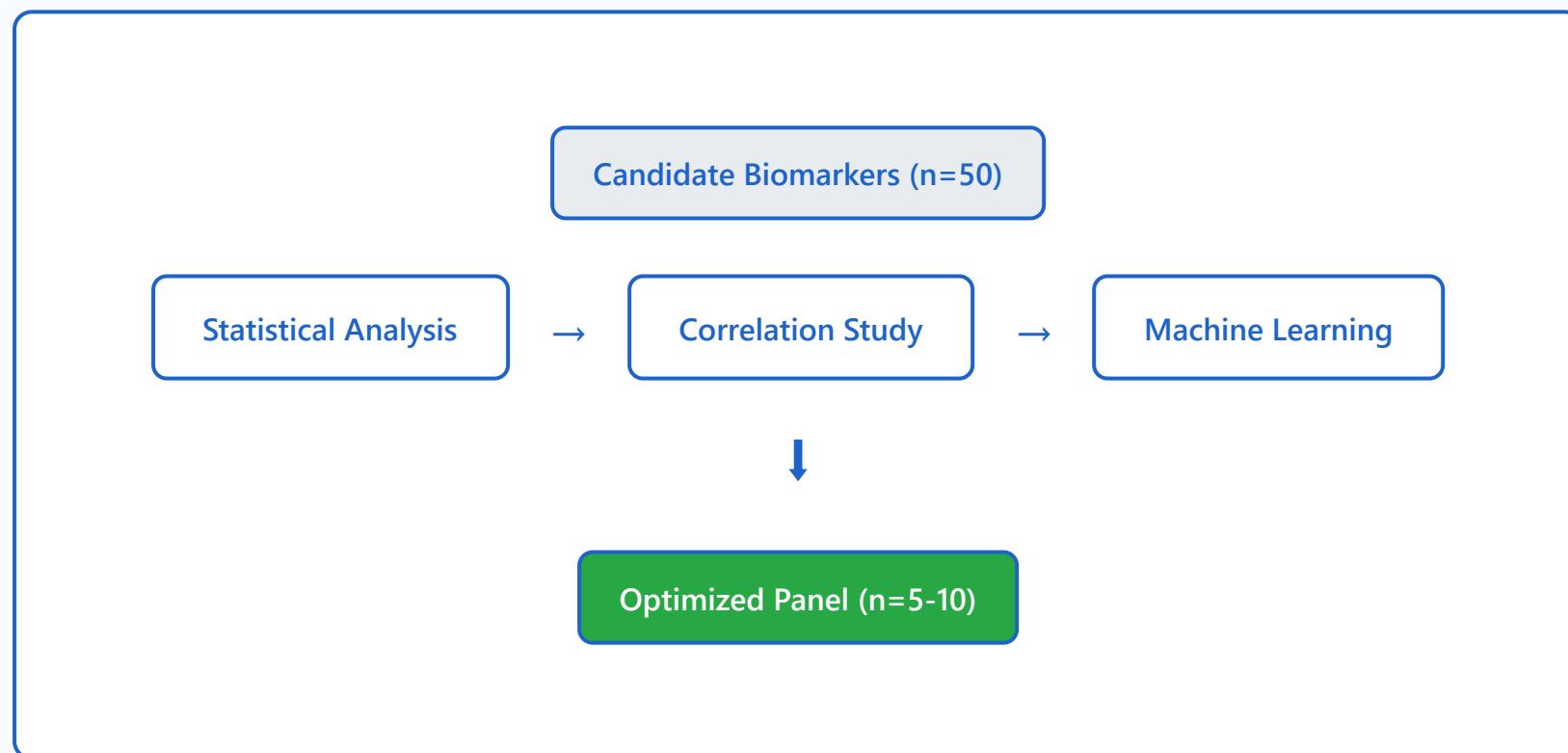
#### Key Advantages:

- Increased sensitivity and specificity compared to single biomarkers
- Ability to detect multiple disease states simultaneously

- Reduced false positive and false negative rates
- Comprehensive assessment of disease mechanisms

## 2. Optimal Combinations

Selecting the optimal combination of biomarkers is crucial for panel development. Feature selection algorithms identify the most informative biomarkers while minimizing redundancy and complexity.



**Selection Methods:**

Method	Approach	Application
Univariate Analysis	Individual marker performance	Initial screening
Multivariate Models	Logistic regression, Cox models	Risk prediction
Machine Learning	Random forests, SVM, Neural networks	Complex pattern recognition
LASSO/Ridge	Penalized regression	Feature reduction

**Example: Cancer Detection Panel**

**Initial candidates:** 50 protein markers

**After correlation analysis:** 20 independent markers

**After ML optimization:** 7-marker panel

**Result:** 92% sensitivity, 95% specificity

### 3. Performance Metrics

Evaluating biomarker panel performance requires multiple metrics to assess clinical utility. These metrics help determine the panel's ability to correctly identify disease presence or absence.

Disease Status

Positive

Negative

Test Result	Positive	True Positive (TP)	False Positive (FP)
	Negative	False Negative (FN)	True Negative (TN)

### Key Performance Metrics:

Metric	Formula	Clinical Meaning	Target
Sensitivity	$TP / (TP + FN)$	Ability to detect disease	>90%
Specificity	$TN / (TN + FP)$	Ability to rule out disease	>90%
PPV	$TP / (TP + FP)$	Probability of disease if test positive	>80%
NPV	$TN / (TN + FN)$	Probability of no disease if test negative	>95%
AUC-ROC	Area under curve	Overall discrimination ability	>0.85

#### Example: Breast Cancer Screening Panel

Study population: 1,000 women (100 with cancer)

Test results: 95 TP, 5 FN, 810 TN, 90 FP

**Sensitivity:** 95% (95/100)

**Specificity:** 90% (810/900)

**PPV:** 51.4% (95/185)

**NPV:** 99.4% (810/815)

## 4. Cost-benefit Analysis

Clinical and economic considerations are essential for successful panel implementation. Cost-effectiveness analysis balances diagnostic performance with healthcare resource utilization.

Test Cost

Reagent Cost

Labor Cost



Total Direct Costs

VS

Early Detection

Prevent Complications

Improved Outcomes



Healthcare Savings

## Economic Evaluation Components:

Category	Cost Factors	Benefit Factors
Direct Medical	Test materials, equipment, personnel	Avoided hospitalizations, reduced treatment costs
Indirect	Follow-up tests, false positive workup	Improved productivity, reduced disability
Intangible	Patient anxiety, waiting time	Quality of life, peace of mind

### Example: Cardiovascular Risk Panel

**Panel cost:** \$200 per test

**Population:** 10,000 high-risk patients

**Total testing cost:** \$2,000,000

**Prevented MI/Stroke:** 150 events

**Average treatment cost saved:** \$50,000 per event

**Total savings:** \$7,500,000

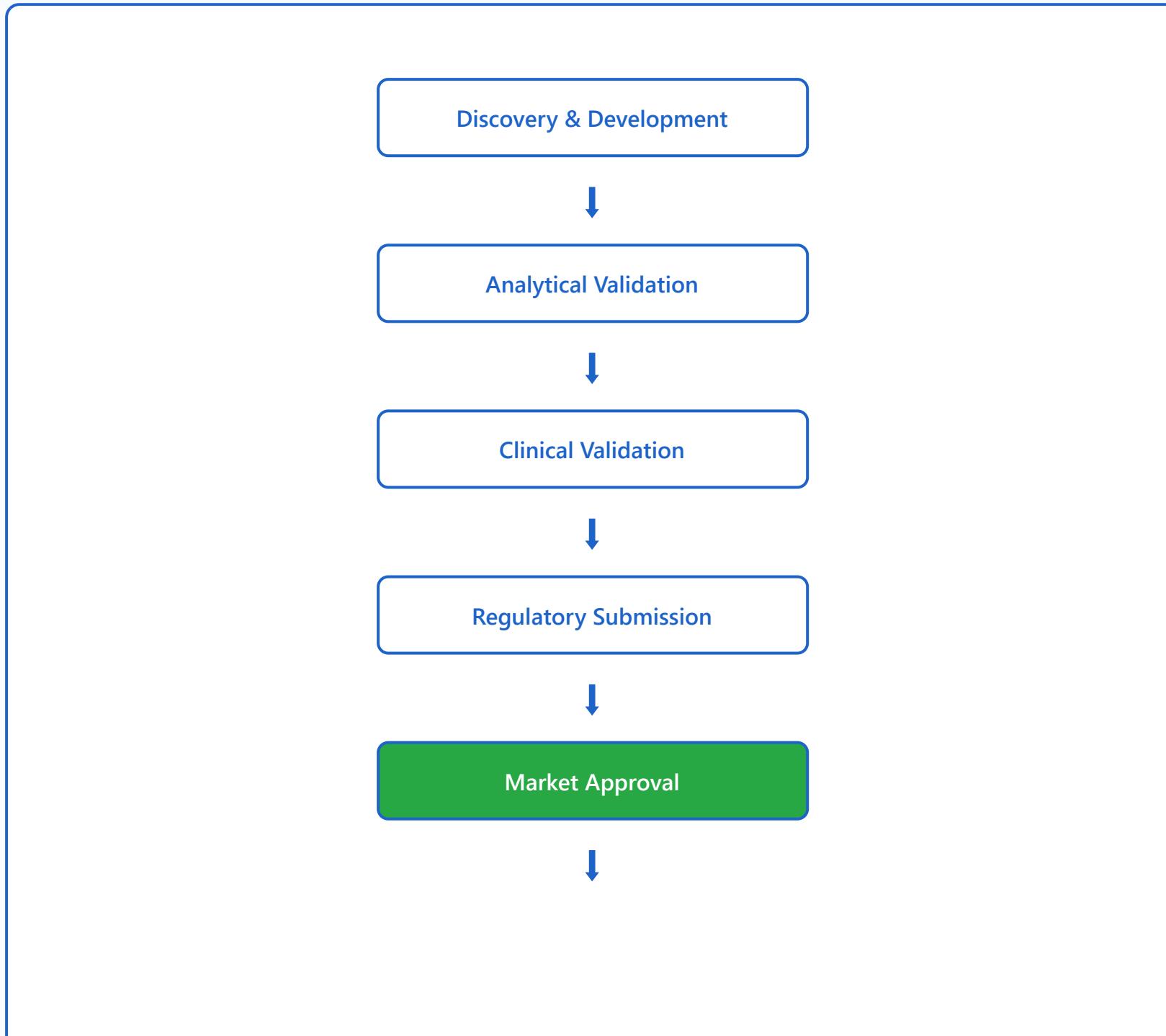
**Net benefit:** \$5,500,000 (ROI: 275%)

## Clinical Considerations:

- Impact on treatment decisions and patient outcomes
- Integration into existing clinical workflows
- Training requirements for healthcare providers
- Patient acceptance and compliance

## 5. Regulatory Approval

Biomarker panels must navigate complex regulatory pathways to reach clinical use. FDA and EMA approval processes ensure safety, efficacy, and clinical validity.



**FDA Approval Pathways:**

Pathway	Description	Timeline	Examples
510(k) Clearance	Substantially equivalent to existing device	3-6 months	Routine diagnostic panels
De Novo	Novel low-to-moderate risk devices	6-12 months	New biomarker combinations
PMA	High-risk devices requiring clinical trials	1-3 years	Cancer screening panels
LDT	Laboratory Developed Tests (CLIA-certified)	Variable	Specialty laboratory tests

**Example: Oncotype DX (Breast Cancer Panel)****Type:** 21-gene expression panel**Pathway:** Initially LDT, later FDA De Novo approval**Clinical validation:** Multiple prospective studies (10,000+ patients)**Evidence required:** Analytical validity, clinical validity, clinical utility**Result:** FDA approval + Medicare coverage**Key Regulatory Requirements:**

- **Analytical Validation:** Accuracy, precision, sensitivity, specificity, reproducibility
- **Clinical Validation:** Association with clinical outcomes in target population
- **Clinical Utility:** Impact on patient management and health outcomes

- **Quality Systems:** Manufacturing, quality control, documentation
- **Labeling:** Intended use, limitations, interpretation guidance

#### **EMA Approval (European Union):**

- In Vitro Diagnostic Regulation (IVDR) - implemented 2022
- Risk classification: Class A (lowest) to Class D (highest)
- CE marking required for market access
- Notified Body assessment for higher-risk devices
- Post-market surveillance and vigilance reporting

# Systems Medicine

A Network-Based Approach to Understanding and Treating Disease

## Network Medicine

Disease as network perturbations

## Disease Modules

Interconnected disease components

## Comorbidities

Shared molecular mechanisms

## Drug Repurposing

Network-based drug discovery

## Personalized Networks

Patient-specific network models

## Detailed Explanations

1

### Network Medicine

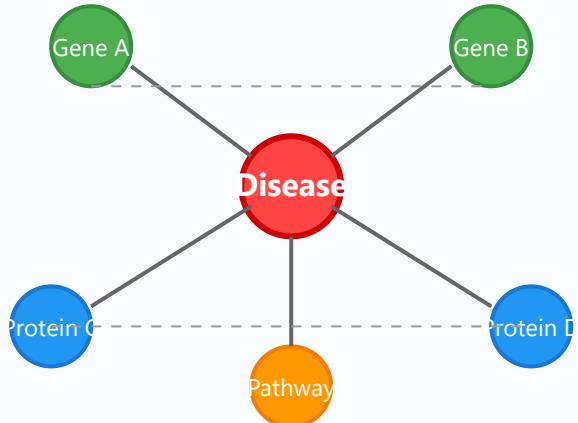
Network Medicine is a paradigm that views diseases not as isolated events, but as perturbations in complex biological networks. This approach integrates genomics, proteomics, and metabolomics data to understand disease mechanisms.

### Key Concepts:

- **Interactome:** The complete set of molecular interactions in cells, including protein-protein, protein-DNA, and metabolic interactions
- **Network Perturbation:** Disease occurs when normal network structure or function is disrupted
- **Systems-level Analysis:** Understanding how multiple components interact simultaneously

### Applications:

- Identifying disease-associated genes and pathways
- Predicting disease progression
- Understanding complex disease etiology



### Key Insight:

Diseases rarely result from abnormalities in a single gene or protein, but from complex interactions within cellular networks.

# Disease Modules

Disease modules are interconnected groups of genes, proteins, or other molecular components that collectively contribute to a disease phenotype. These modules represent functional units within the larger biological network.

## Characteristics:

- **Modularity:** Disease components cluster together in the network
- **Functional Coherence:** Module members share biological functions
- **Topological Proximity:** Disease genes tend to interact directly or through few intermediates

## Significance:

- Reveals disease mechanisms at the systems level
- Identifies potential therapeutic targets
- Explains genetic heterogeneity in diseases
- Enables disease classification based on molecular profiles



## Key Insight:

Disease modules help explain why different genetic mutations can lead to similar disease phenotypes and why diseases share common molecular mechanisms.

### 3 Comorbidities

Comorbidities refer to the co-occurrence of two or more diseases in the same individual more frequently than expected by chance. Network medicine reveals that comorbidities often result from shared molecular mechanisms and overlapping disease modules.

#### Network-Based Explanations:

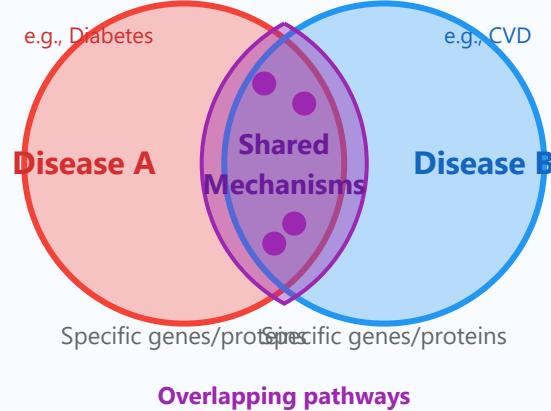
- **Overlapping Disease Modules:** Two diseases share common genes or proteins
- **Pathway Convergence:** Different diseases affect the same biological pathways
- **Shared Risk Factors:** Common genetic or environmental factors

#### Clinical Examples:

- Diabetes and cardiovascular disease
- Depression and chronic pain
- Obesity and metabolic syndrome
- Autoimmune diseases (e.g., rheumatoid arthritis and psoriasis)

#### Implications:

- Improved disease prediction and prevention



- Comprehensive treatment strategies
- Understanding disease progression patterns

### Key Insight:

Network analysis reveals that comorbidities are not random but follow predictable patterns based on shared molecular mechanisms, enabling better patient stratification and treatment planning.

## 4 Drug Repurposing

Network-based drug repurposing identifies new therapeutic uses for existing drugs by analyzing their effects on biological networks. This approach can significantly reduce the time and cost of drug development.

### Network-Based Strategies:

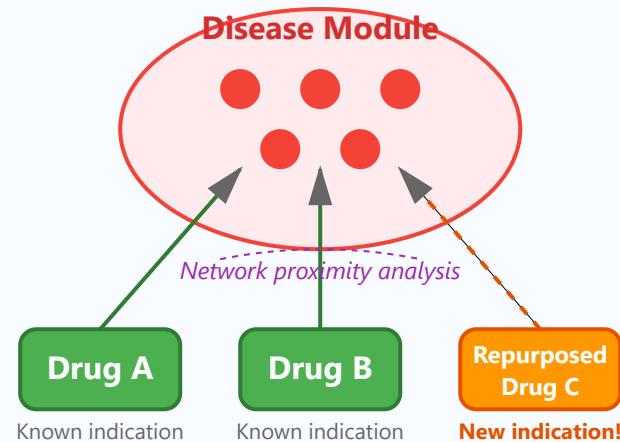
- **Target-Based:** Identifying drugs that target proteins in disease modules
- **Pathway-Based:** Finding drugs that modulate disease-relevant pathways
- **Network Proximity:** Measuring closeness between drug targets and disease genes
- **Module-Based:** Drugs affecting the same network modules as diseases

## Success Stories:

- Aspirin: Originally for pain → cardiovascular protection
- Metformin: Diabetes → cancer prevention, aging
- Thalidomide: Sedative → multiple myeloma treatment
- Sildenafil (Viagra): Hypertension → erectile dysfunction

## Advantages:

- Reduced development time (3-12 years vs 10-17 years)
- Lower costs (\$300M vs \$2-3B)
- Known safety profiles
- Higher success rates



### Key Insight:

By analyzing network proximity between drug targets and disease genes, we can systematically identify repurposing opportunities that would be missed by traditional approaches.

## 5 Personalized Networks

Personalized networks integrate individual patient data (genomics, transcriptomics, proteomics, clinical data) to create patient-specific network models. This enables precision

medicine approaches tailored to each patient's unique molecular profile.

### Data Integration:

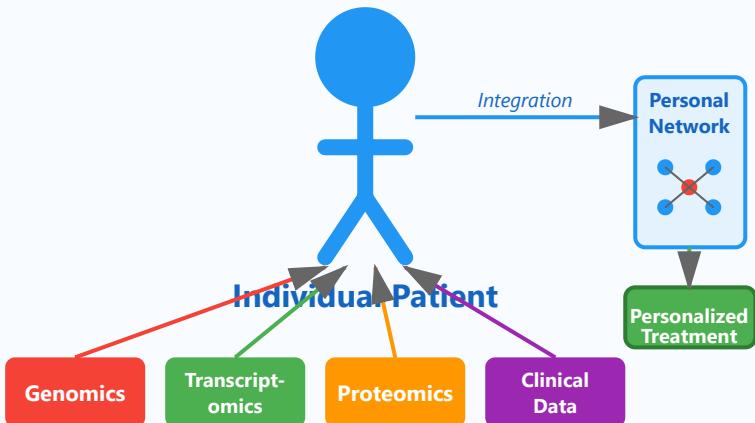
- **Genomic Data:** Mutations, SNPs, structural variants
- **Transcriptomic Data:** Gene expression patterns
- **Proteomic Data:** Protein abundance and modifications
- **Metabolomic Data:** Metabolite profiles
- **Clinical Data:** Symptoms, disease history, treatment response

### Applications:

- Predicting individual disease risk
- Selecting optimal treatments
- Identifying personalized drug targets
- Predicting treatment response and adverse effects
- Monitoring disease progression

### Challenges:

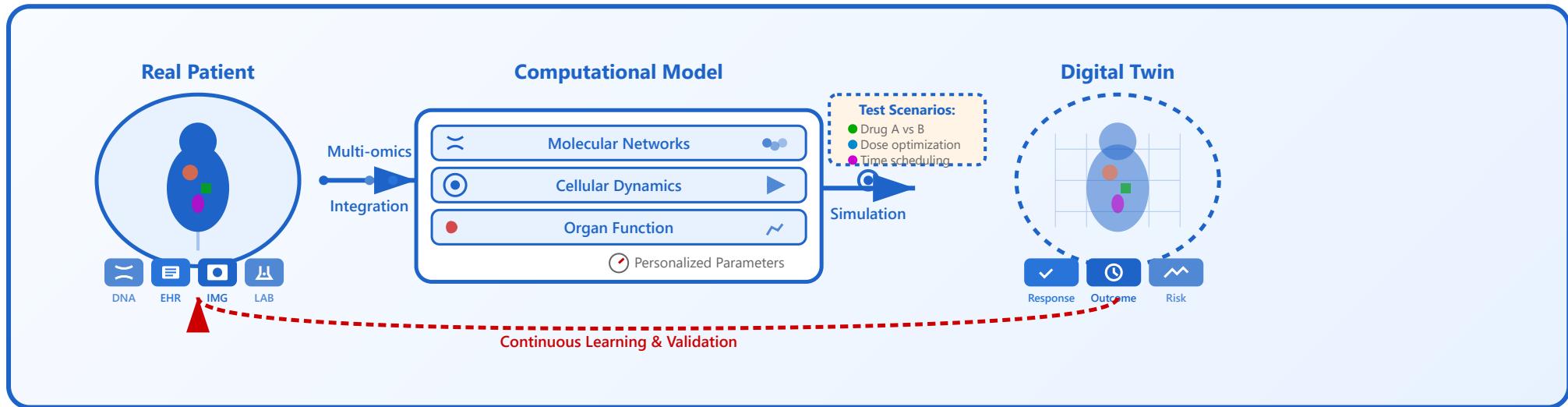
- Data integration complexity
- Computational requirements
- Clinical validation
- Cost and accessibility



**Key Insight:**

Personalized networks enable true precision medicine by capturing each patient's unique molecular landscape and identifying optimal therapeutic strategies based on individual network perturbations.

# Digital Twins in Medicine



## 1. Patient Models

Computational patient representations across multiple biological scales

## 2. Simulation Frameworks

In silico clinical trials for virtual drug testing

## 3. Parameter Estimation

Personalizing model parameters from patient data

## 4. Treatment Optimization

## 5. Validation Approaches

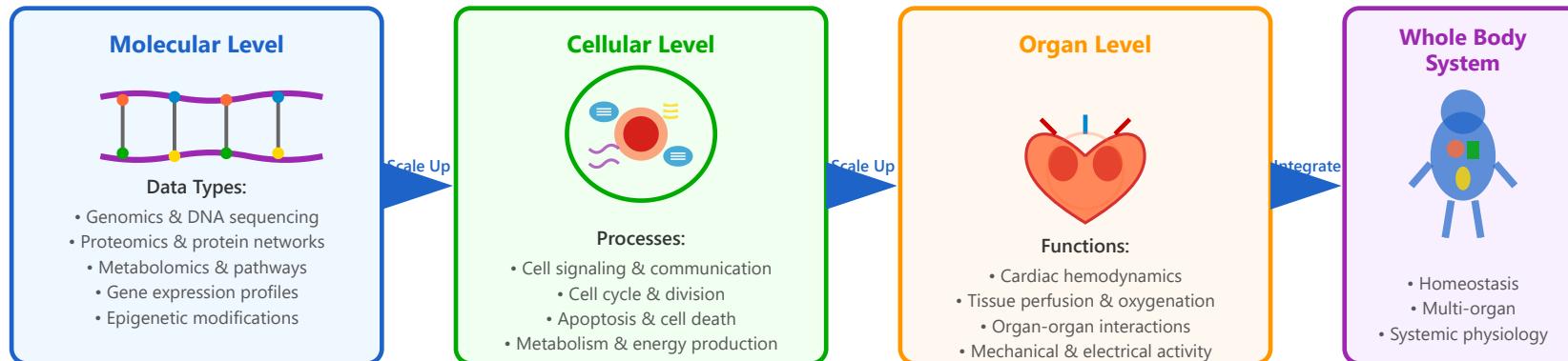
Simulating and optimizing treatment strategies

Comparing predictions to real-world outcomes

# 1

## Patient Models: Computational Patient Representations

### Multi-Scale Patient Modeling Framework



Patient models form the foundation of digital twin technology by creating computational representations that capture the biological complexity of individual patients across multiple scales. These models integrate diverse data

sources including genomics, proteomics, imaging, and clinical parameters to construct a comprehensive digital representation of patient physiology and pathology. The multi-scale approach enables understanding of how molecular changes propagate through cellular, tissue, and organ levels to manifest as clinical symptoms.

## Key Components

- ▶ **Multi-scale Integration:** Models span from molecular interactions at the gene and protein level through cellular processes and organ function to whole-body physiology, capturing the hierarchical nature of biological systems
- ▶ **Mathematical Frameworks:** Utilize ordinary differential equations (ODEs) for temporal dynamics, partial differential equations (PDEs) for spatial distributions, agent-based models for cellular interactions, and machine learning for pattern recognition
- ▶ **Data Assimilation:** Continuously integrate patient-specific data including electronic health records, laboratory results, medical imaging (CT, MRI, PET), genomic sequencing, and wearable sensor data
- ▶ **Dynamic Representation:** Capture temporal evolution of disease states, treatment responses, and physiological changes over multiple time scales from seconds (heartbeat) to years (aging, disease progression)
- ▶ **Uncertainty Quantification:** Explicitly model biological variability and measurement uncertainty to provide confidence intervals on predictions

## Clinical Example: Cardiovascular Digital Twin

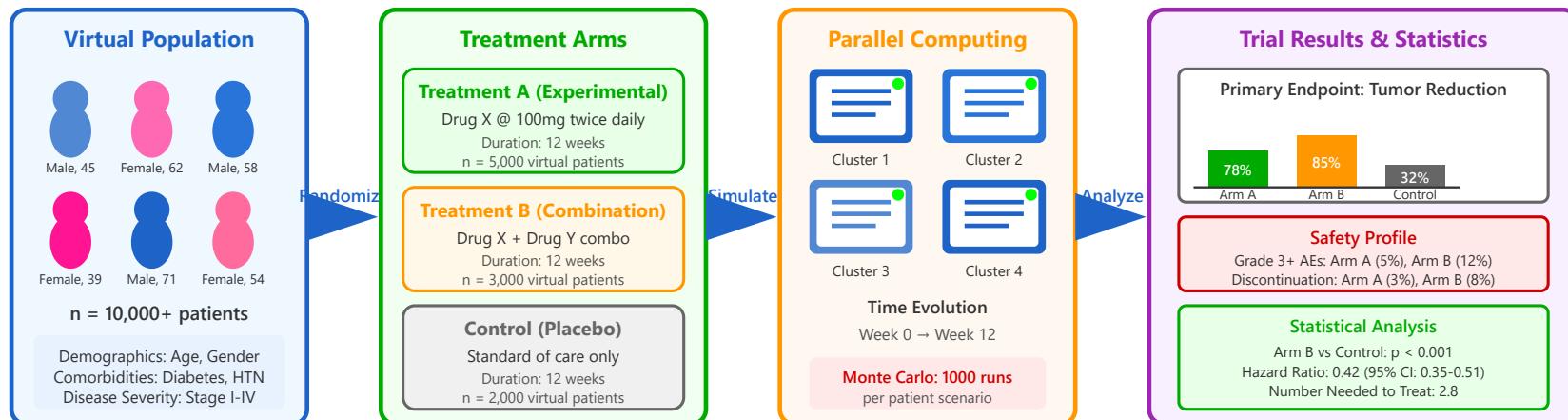
A cardiac digital twin for a 65-year-old patient with heart failure combines CT/MRI imaging to create a patient-specific 3D heart geometry showing chamber dimensions and wall thickness. ECG data calibrates electrical conduction parameters, while blood biomarkers (BNP, troponin) inform metabolic state. Genetic variants in ion channel genes (SCN5A, KCNH2) personalize electrophysiology parameters. The model successfully predicts ejection fraction changes with beta-blockers (predicted: 38% → 45%, observed: 37% → 46%), simulates stress test responses to guide safe exercise limits, and identifies optimal cardiac resynchronization therapy pacing parameters, reducing hospitalizations by 60% over 12 months.

## 2

## Simulation Frameworks: In Silico Clinical Trials

---

## Virtual Clinical Trial Pipeline



Computational time: 48 hours vs 2-3 years for traditional clinical trial

In silico clinical trials leverage computational simulation frameworks to conduct virtual studies that complement or reduce the need for traditional clinical trials. These frameworks enable researchers to test hypotheses, optimize trial designs, identify optimal patient populations, and predict treatment outcomes before conducting expensive and time-consuming real-world trials. By simulating thousands of virtual patients, researchers can rapidly explore multiple treatment strategies and identify the most promising approaches for further development.

### Key Advantages

- ▶ **Accelerated Development:** Rapidly test multiple treatment strategies, doses, combinations, and schedules simultaneously without patient safety concerns, reducing development time from years to weeks

- ▶ **Cost Reduction:** Dramatically reduce trial costs (up to 90%) by identifying failed approaches early and optimizing trial parameters before real patient enrollment, saving millions in development costs
- ▶ **Ethical Benefits:** Minimize patient exposure to potentially ineffective or harmful treatments through virtual screening, protecting vulnerable populations from unnecessary risk
- ▶ **Rare Disease Research:** Enable trials for conditions where patient recruitment is challenging by generating synthetic cohorts that match real population statistics
- ▶ **Regulatory Insights:** Provide supporting evidence for regulatory submissions (FDA, EMA) and help design more informative and efficient clinical trials
- ▶ **Dose Optimization:** Identify optimal dosing regimens and treatment schedules before human trials, maximizing efficacy while minimizing toxicity

## Clinical Example: Oncology Drug Development

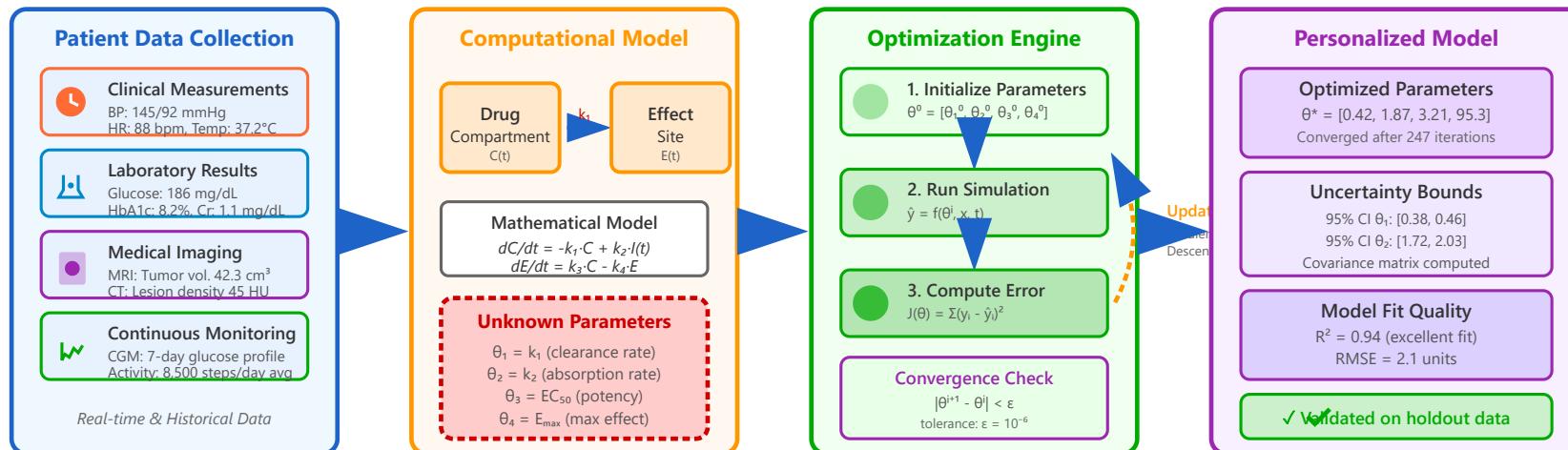
A pharmaceutical company developing a novel cancer immunotherapy uses in silico trials to test the drug across 10,000 virtual patients with varying tumor genetics (KRAS, BRAF, TP53 mutations), immune profiles (PD-L1 expression, tumor infiltrating lymphocytes), and disease stages (I-IV). The simulation identifies that the drug is most effective in patients with high PD-L1 expression (>50%) combined with specific genetic mutations. This enables the company to design a biomarker-enriched Phase II trial enrolling only likely responders, reducing trial size from 500 to 150 patients while increasing success probability from 30% to 65%. The virtual trial also identifies an optimal dosing schedule (10mg/kg).

Q3W) that balances efficacy and immune-related adverse events, and predicts a median overall survival benefit of 8.4 months vs control, which is subsequently confirmed in the real trial (observed: 8.1 months, 95% CI: 6.8-9.7).

### 3

## Parameter Estimation: Personalizing Model Parameters

### Parameter Estimation & Model Calibration Workflow



Parameter estimation is the critical process of calibrating computational models to match individual patient characteristics and responses. This personalization transforms generic models into patient-specific digital twins by determining the optimal parameter values that best reproduce observed clinical data. The process involves sophisticated mathematical optimization and statistical methods including maximum likelihood estimation, Bayesian inference, and nonlinear least squares to ensure model accuracy and quantify uncertainty in parameter estimates.

## Key Methods & Techniques

- ▶ **Maximum Likelihood Estimation (MLE):** Finds parameters that maximize the probability of observing the measured data given the model, assuming a specific error distribution (typically Gaussian)
- ▶ **Bayesian Inference:** Incorporates prior knowledge about parameter ranges and quantifies parameter uncertainty through posterior probability distributions using Markov Chain Monte Carlo (MCMC) sampling
- ▶ **Optimization Algorithms:** Employs gradient descent (Levenberg-Marquardt), genetic algorithms, particle swarm optimization, simulated annealing, or ensemble Kalman filters to efficiently search high-dimensional parameter spaces
- ▶ **Identifiability Analysis:** Determines which parameters can be reliably estimated from available data using sensitivity analysis, Fisher Information Matrix, and profile likelihood methods
- ▶ **Sensitivity Analysis:** Assesses how parameter uncertainty propagates to model predictions, guiding data collection priorities and identifying which measurements are most informative
- ▶ **Cross-Validation:** Uses k-fold or leave-one-out cross-validation to prevent overfitting and ensure parameter

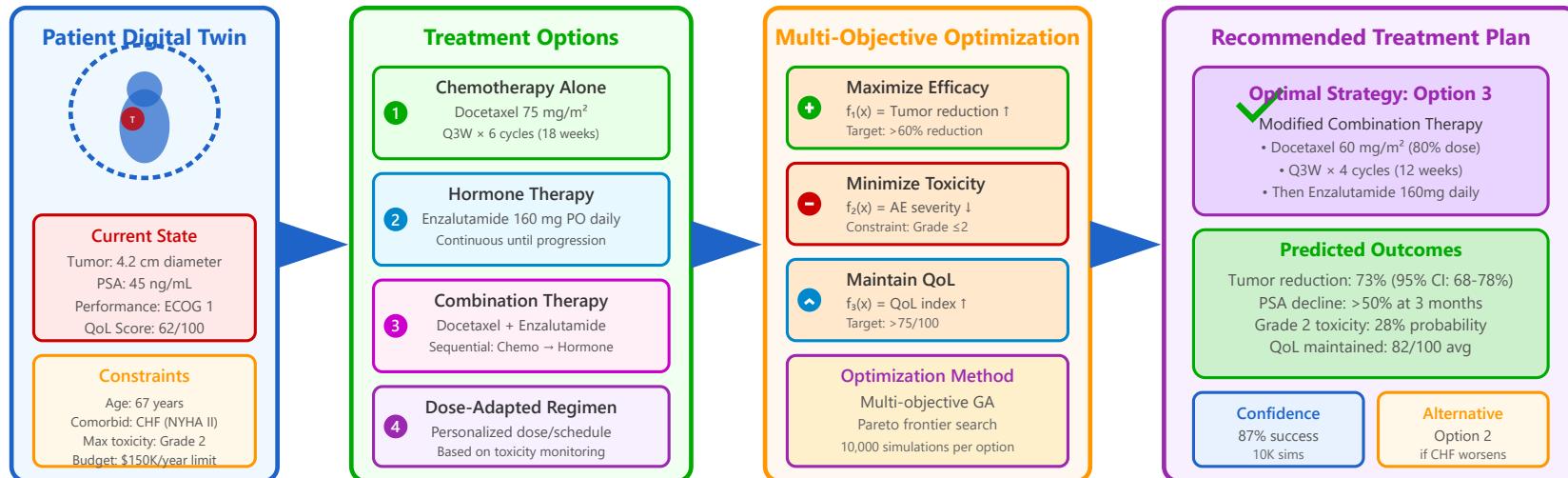
## Clinical Example: Diabetes Management Personalization

For a 54-year-old Type 2 diabetes patient with BMI  $32 \text{ kg/m}^2$ , parameter estimation uses 14 days of continuous glucose monitoring data (CGM), detailed insulin dosing history (long-acting insulin glargine 28 units daily, rapid-acting insulin aspart with meals), meal logs with carbohydrate counting, and exercise records from wearable devices. The optimization algorithm calibrates a glucose-insulin dynamics model, determining patient-specific parameters: insulin sensitivity ( $\text{SI} = 0.00042 \text{ L/mU/min}$ , 30% below population average), glucose effectiveness ( $\text{SG} = 0.021 \text{ min}^{-1}$ ), basal glucose production ( $\text{EGP}_0 = 2.1 \text{ mg/kg/min}$ ), and insulin clearance rate ( $k = 0.15 \text{ min}^{-1}$ ). These calibrated parameters enable accurate prediction of glucose responses to meals and insulin with RMSE of 12 mg/dL. The personalized model supports optimized dosing strategies including adjusted carbohydrate ratios (1:12g instead of standard 1:15g) and correction factors, reducing hypoglycemic events below 70 mg/dL by 45% (from 3.2 to 1.8 events/week) and improving time-in-range 70-180 mg/dL from 62% to 78% over 3 months.

4

## Treatment Optimization: Simulating Treatment Strategies

## Multi-Objective Treatment Optimization Framework



Treatment optimization leverages digital twins to identify the best therapeutic strategy for individual patients by simulating and comparing multiple treatment options across multiple competing objectives. This process goes beyond simple efficacy evaluation to balance effectiveness, safety, quality of life, treatment duration, cost, and patient preferences. The optimization employs advanced algorithms including multi-objective genetic algorithms, Pareto optimization, and reinforcement learning to explore vast treatment spaces and identify optimal strategies that account for patient-specific constraints, comorbidities, and individual goals.

### Key Capabilities & Methods

- ▶ **Multi-Objective Optimization:** Simultaneously optimizes multiple competing goals (efficacy, safety, QoL, cost) using Pareto frontier analysis to identify treatments offering best trade-offs across all objectives
- ▶ **Dose and Schedule Optimization:** Determines optimal drug dosing regimens, radiation fractionation schedules, timing of combination therapies, and treatment holidays using pharmacokinetic/pharmacodynamic modeling to maximize therapeutic window
- ▶ **Adaptive Planning:** Enables mid-treatment adjustments based on observed responses through Bayesian updating, allowing personalized adaptation as patient condition evolves (e.g., dose reduction for toxicity, dose intensification for poor response)
- ▶ **Risk Stratification:** Quantifies uncertainty in outcome predictions using Monte Carlo simulation and sensitivity analysis, identifying patient-specific risk factors that may affect treatment success
- ▶ **Alternative Scenarios:** Explores "what-if" scenarios (disease progression, treatment complications, comorbidity changes) to understand robustness and prepare contingency plans
- ▶ **Constraint Handling:** Incorporates hard constraints (maximum toxicity, budget limits) and soft preferences (minimize treatment visits, oral vs IV preference) into optimization

## Clinical Example: Precision Radiation Oncology for Non-Small Cell Lung Cancer

For a 68-year-old lung cancer patient with a 4.8 cm right upper lobe tumor adjacent to critical structures, the digital twin integrates 4D-CT imaging showing tumor motion with respiration (1.2 cm superior-inferior), PET imaging indicating heterogeneous FDG uptake (SUVmax 12.8), pulmonary function tests (FEV1 68% predicted), and cardiac assessment showing prior MI with EF 45%. The optimization algorithm explores 50,000 different radiation delivery plans varying:

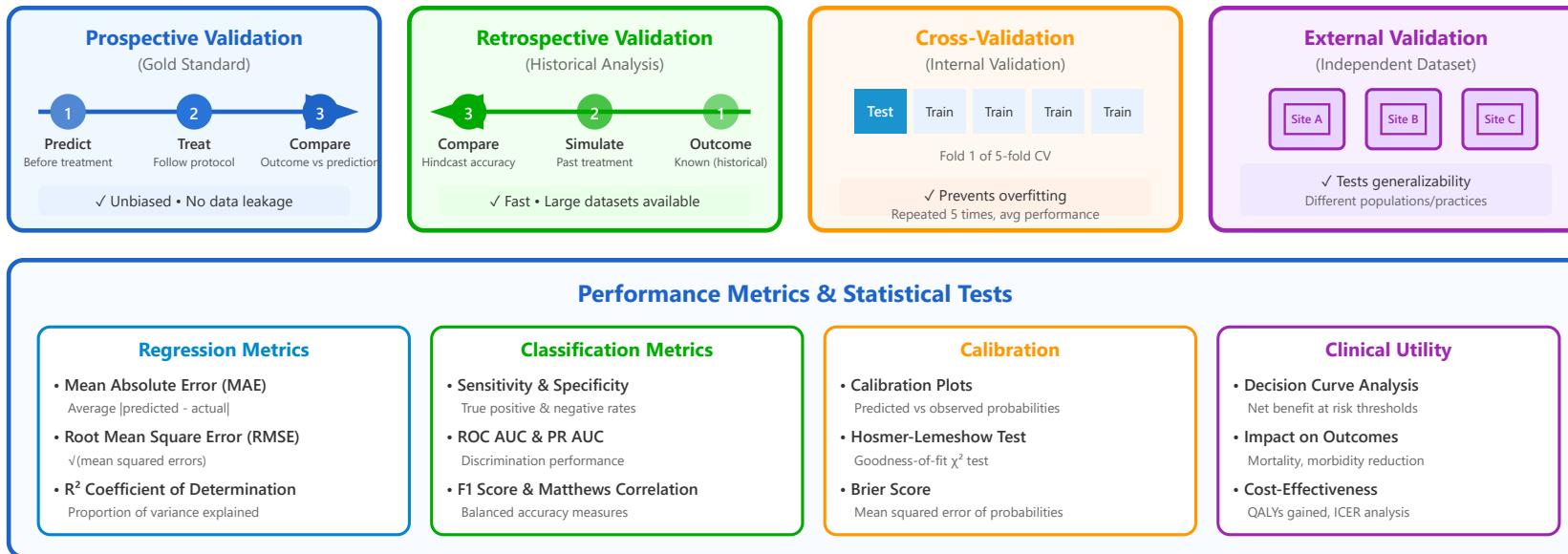
beam angles (7-11 coplanar/non-coplanar), intensities (IMRT vs VMAT), fractionation (60Gy/30fx vs 70Gy/35fx vs SBRT 50Gy/5fx), and motion management strategies (free breathing vs breath-hold vs gating). The optimal plan delivers 70 Gy in 35 fractions using dual-arc VMAT with breath-hold, achieving: tumor D95% = 68 Gy (97% coverage), heart V30 = 18% (goal <30%), spinal cord Dmax = 42 Gy (limit <45 Gy), mean lung dose = 14 Gy (limit <20 Gy). Simulations predict 85% local control at 2 years (vs 78% with standard planning), 8% risk of grade 3+ pneumonitis (vs 23% standard), and 4% cardiac event risk at 5 years (vs 9% standard). The plan is successfully delivered with observed pneumonitis rate of 7.5% and 2-year local control of 83%, validating the digital twin predictions.

## 5

## Validation Approaches: Comparing Predictions to Reality

---

# Comprehensive Validation Framework



Validation is the cornerstone of digital twin credibility, providing essential evidence that computational predictions reliably match real-world patient outcomes. Rigorous validation involves multiple complementary approaches including prospective testing on new patients (the gold standard), retrospective analysis of historical data (for rapid assessment), cross-validation within datasets (to prevent overfitting), and external validation across different institutions and populations (to ensure generalizability). Without thorough validation using multiple independent datasets and diverse patient populations, digital twins remain theoretical constructs rather than trusted clinical decision support tools worthy of regulatory approval and clinical adoption.

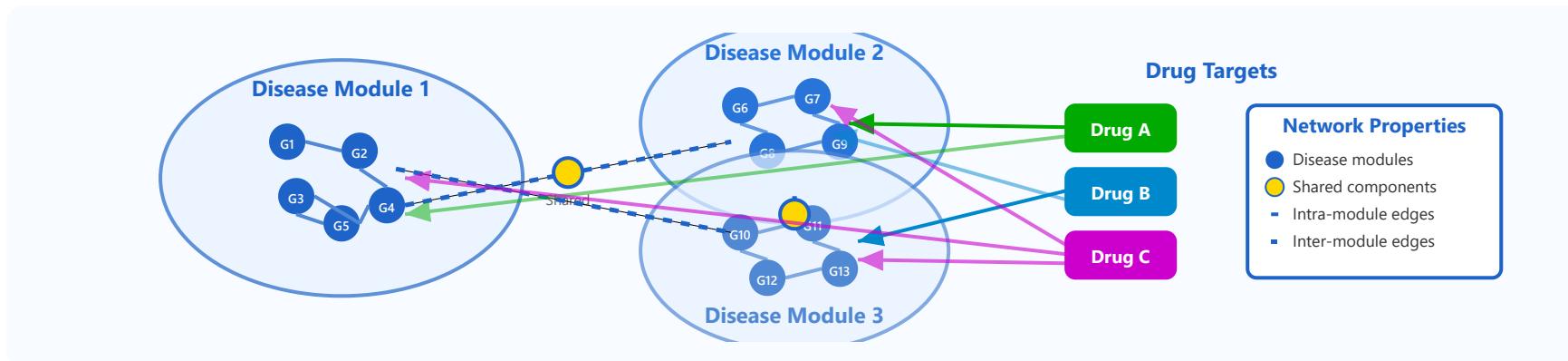
## Key Validation Strategies

- ▶ **Prospective Clinical Validation:** The gold standard where predictions are made before treatment and subsequently compared to actual observed outcomes in real patients receiving care. This approach eliminates data leakage and provides unbiased assessment of real-world performance
- ▶ **Retrospective Analysis:** Tests model performance on historical patient data where outcomes are already known, useful for initial validation and model refinement before expensive prospective testing. Enables rapid iteration and identification of systematic biases
- ▶ **Cross-Validation:** Systematically partitions data into training and testing sets (k-fold, leave-one-out, stratified sampling) to assess model generalization and prevent overfitting to specific patient cohorts or center-specific practices
- ▶ **External Validation:** Tests models on completely independent datasets from different institutions, geographic regions, time periods, or populations to ensure robustness and broad applicability beyond the development cohort
- ▶ **Temporal Validation:** Evaluates model performance on patients treated in different time periods to assess stability as medical practice evolves and treatment standards change
- ▶ **Subgroup Analysis:** Examines performance across clinically relevant subgroups (age, sex, disease stage, comorbidities) to identify populations where the model performs well or poorly
- ▶ **Sensitivity Analysis:** Evaluates how prediction accuracy varies with model assumptions, parameter uncertainty, and data quality to identify model vulnerabilities and required data quality thresholds

## Clinical Example: Sepsis Prediction System Multi-Phase Validation

A digital twin for early sepsis prediction in ICU patients undergoes comprehensive multi-phase validation demonstrating clinical utility. **Phase 1 (Retrospective):** Analysis on 5,000 ICU admissions from a single center (2018-2020) shows 85% sensitivity and 92% specificity for predicting sepsis onset 6 hours before clinical recognition (SOFA score  $\geq 2$ ), with AUROC 0.91 (95% CI: 0.89-0.93) and calibration slope 1.02. **Phase 2 (Cross-Validation):** 5-fold stratified cross-validation confirms consistent performance across patient subgroups: medical ICU (AUROC 0.89), surgical ICU (0.92), age  $< 65$  (0.90), age  $\geq 65$  (0.91), APACHE II  $< 15$  (0.88), APACHE II  $\geq 15$  (0.93). **Phase 3 (Prospective):** Prospective validation in 500-patient cohort at same institution (2021) achieves 82% sensitivity, 88% specificity, with 15% false alarm rate (median 2.3 alarms/patient), and median advance warning time of 5.8 hours. **Phase 4 (External Validation):** Multi-center validation at three independent hospitals demonstrates maintained performance: Site A (academic, n=800): AUROC 0.87, Site B (community, n=650): AUROC 0.89, Site C (international, n=420): AUROC 0.91. **Phase 5 (Clinical Utility):** Randomized controlled trial (n=1,200, 600 per arm) shows the system enables 3.2 hours earlier treatment initiation (antibiotics, fluids, vasopressors), reducing ICU mortality from 18.3% to 13.1% ( $p=0.002$ , NNT=19), decreasing ICU length of stay from 8.4 to 7.1 days ( $p=0.01$ ), and estimated cost savings of \$4,200 per patient. The comprehensive validation across multiple phases, institutions, and patient populations, combined with demonstrated clinical impact in an RCT, provides strong evidence supporting clinical adoption and regulatory approval.

# Network Medicine



## Disease Networks

Molecular interaction networks in disease

## Interactome

Protein-protein interaction networks

## Disease-disease Relationships

Shared pathways and comorbidities

## Drug-target Networks

Polypharmacology and off-targets

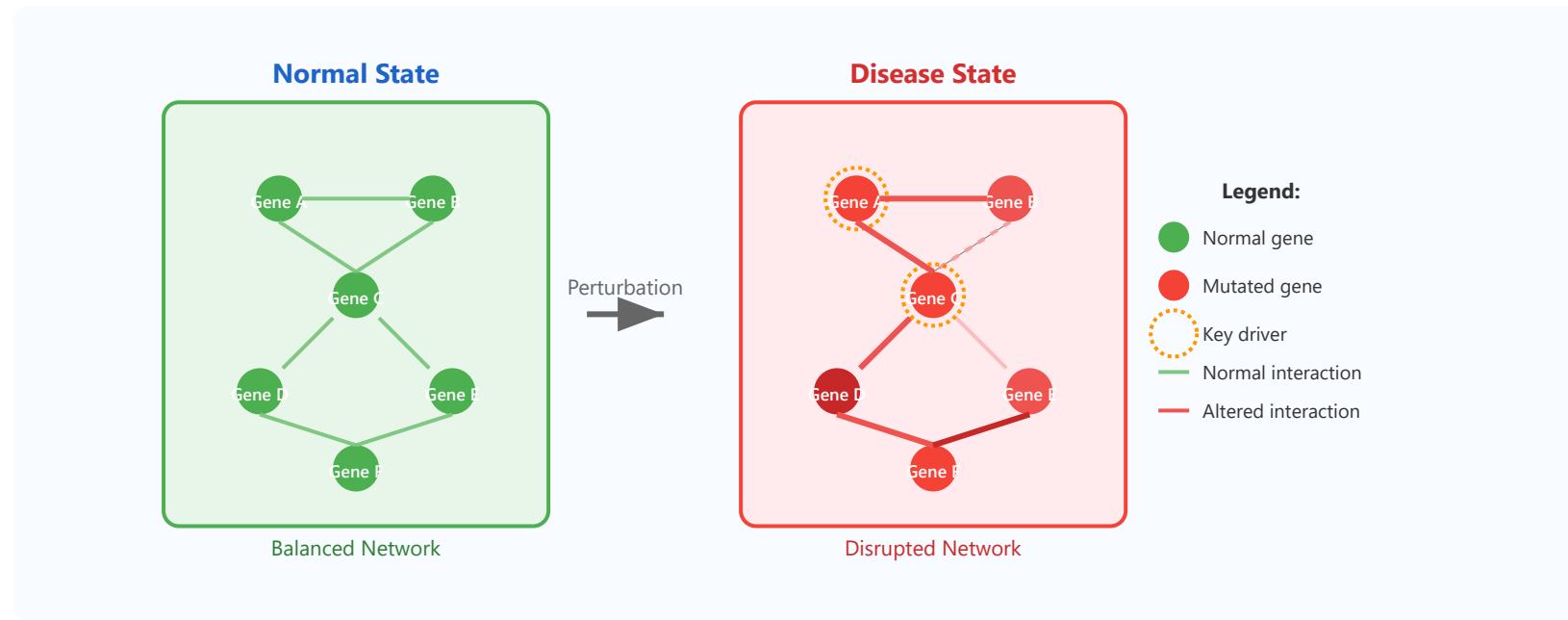
## Network Pharmacology

Systems-level drug discovery

## Detailed Concepts in Network Medicine

## • 1. Disease Networks

Disease networks represent the complex molecular interactions underlying disease pathology. Rather than viewing diseases as caused by single gene defects, network medicine recognizes that most diseases arise from perturbations in interconnected molecular networks. These networks include genes, proteins, metabolites, and their regulatory relationships that collectively contribute to disease phenotypes.



### Key Concepts:

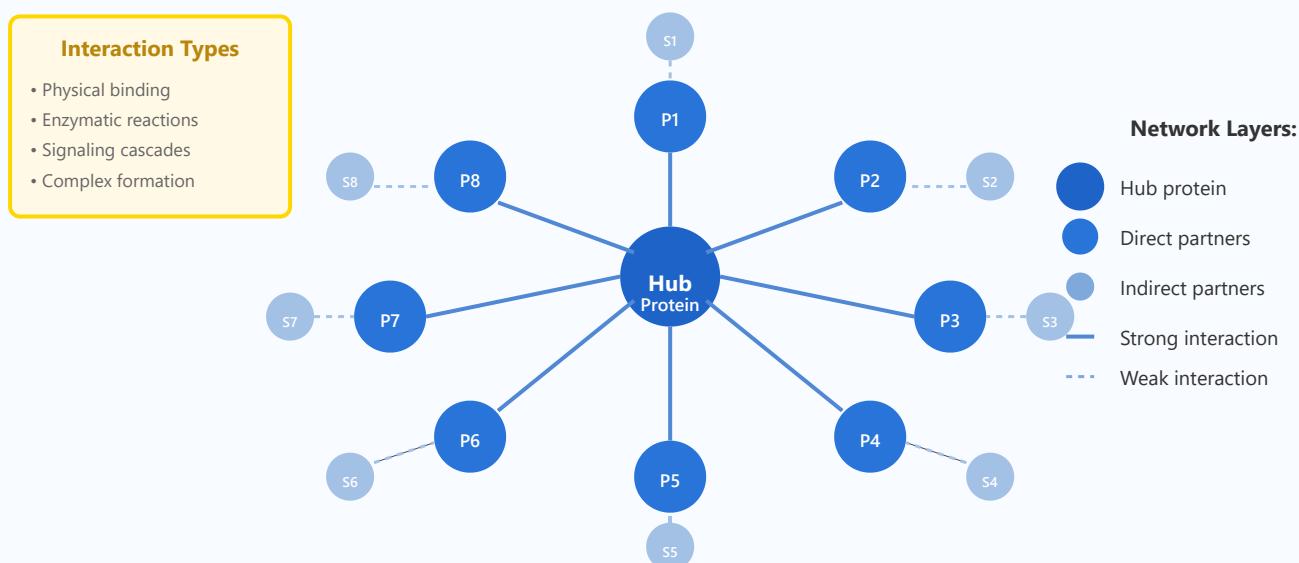
- ▶ Disease modules: Groups of genes and proteins that work together and are associated with specific diseases
- ▶ Network perturbation: Changes in network structure or dynamics that lead to disease phenotypes
- ▶ Pathway analysis: Identifying enriched biological pathways within disease-associated gene sets
- ▶ Hub genes: Highly connected nodes that often play critical roles in disease progression

### Clinical Examples:

Cancer networks reveal how mutations in driver genes like TP53, KRAS, and EGFR disrupt cell cycle regulation and apoptosis pathways. In Alzheimer's disease, networks show how amyloid-beta production, tau phosphorylation, and neuroinflammation are interconnected processes rather than isolated events.

## • 2. Interactome

The interactome represents the comprehensive map of all protein-protein interactions (PPIs) within a cell or organism. These physical interactions form the functional backbone of cellular processes, and understanding them is crucial for deciphering how cells operate in health and disease. The human interactome is estimated to contain over 130,000 protein-protein interactions.



### Key Concepts:

- ▶ Scale-free topology: Few highly connected hubs and many low-degree nodes characterize biological networks
- ▶ Protein complexes: Stable multi-protein assemblies that perform specific cellular functions

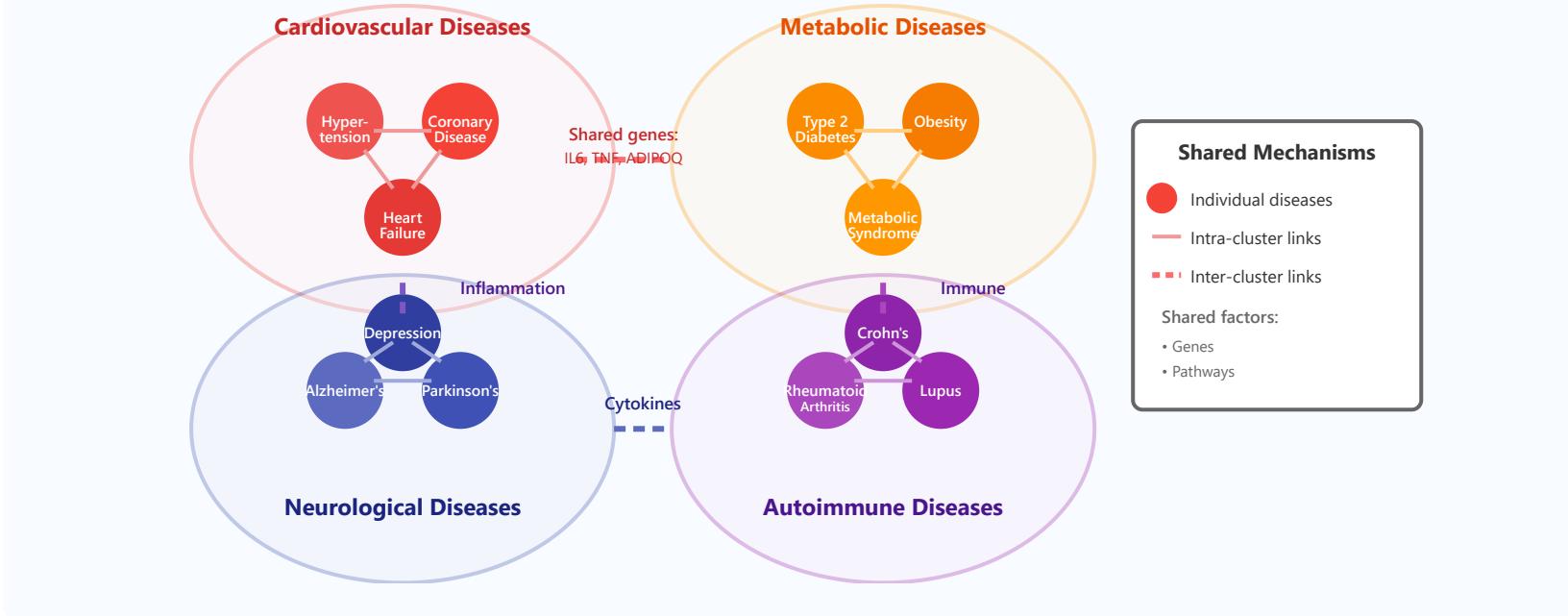
- ▶ Transient interactions: Dynamic protein associations regulated by cellular conditions and signals
- ▶ Network motifs: Recurring patterns of interactions that perform specific information processing tasks

### Clinical Examples:

The p53 protein serves as a major hub in the tumor suppressor network, interacting with over 300 partners. Disruption of key interactome components leads to various diseases such as viral infections hijacking host protein interactions, or mutations in scaffold proteins causing developmental disorders.

### • 3. Disease-disease Relationships

Diseases do not exist in isolation but often share molecular mechanisms, genetic factors, and environmental risks. Understanding these relationships through network analysis reveals patterns of comorbidity, helps identify disease subtypes, and uncovers opportunities for drug repurposing. The human disease network connects diseases through shared genes, proteins, metabolites, or environmental factors.



### Key Concepts:

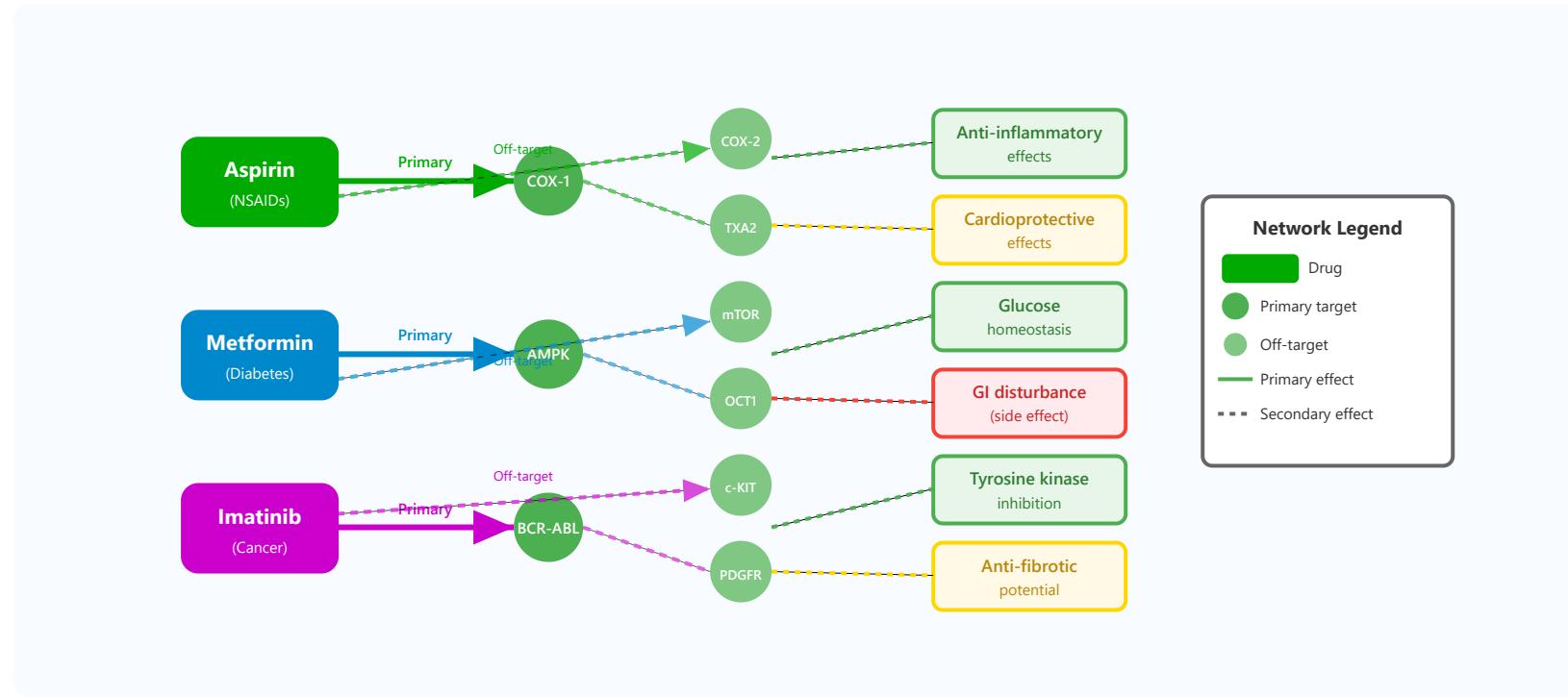
- ▶ Comorbidity patterns: Statistical associations between diseases that co-occur more frequently than expected by chance
- ▶ Disease modules: Overlapping molecular components between diseases indicate shared pathophysiology
- ▶ Pleiotropy: Single genes affecting multiple disease phenotypes reveal fundamental biological connections
- ▶ Disease trajectories: Sequential patterns of disease development that can be predicted from network topology

### Clinical Examples:

Type 2 diabetes and cardiovascular disease share inflammation-related genes and metabolic pathways, explaining their high comorbidity. Inflammatory bowel disease and rheumatoid arthritis share immune dysregulation mechanisms, leading to successful drug repurposing between these conditions.

## • 4. Drug-target Networks

Most drugs interact with multiple proteins beyond their intended targets, creating complex drug-target networks. Understanding these interactions is crucial for predicting drug efficacy, side effects, and discovering new therapeutic applications. Network analysis reveals that successful drugs often target multiple components within disease modules or bridge between related disease pathways.



### Key Concepts:

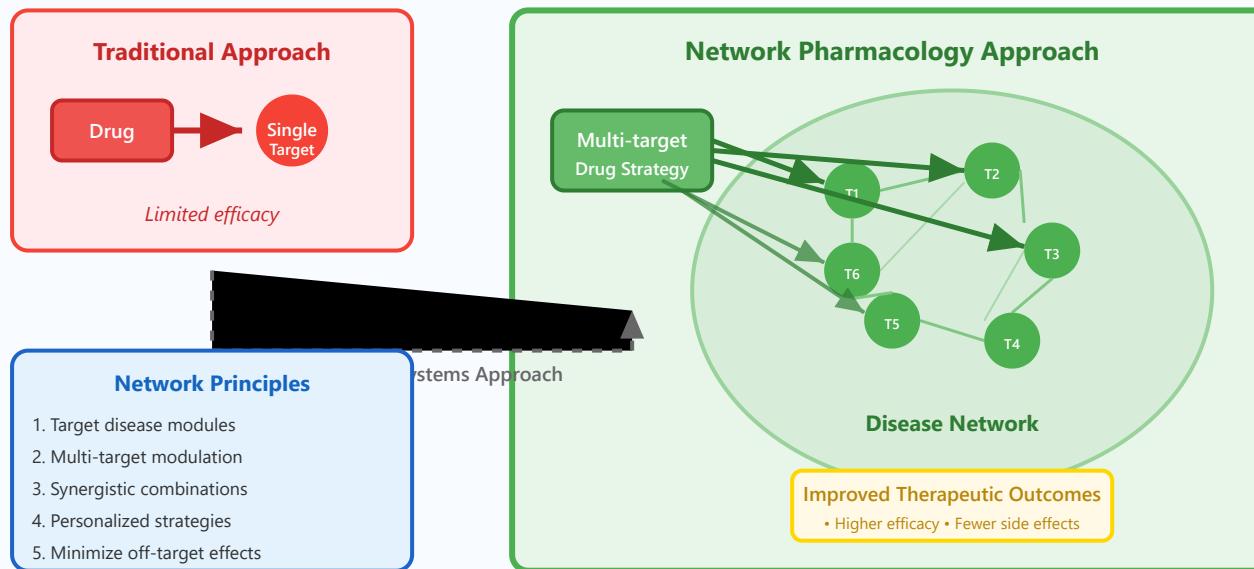
- ▶ Polypharmacology: The ability of drugs to modulate multiple targets, which can enhance efficacy or cause side effects
- ▶ Drug promiscuity: Off-target binding patterns that can be exploited for drug repurposing
- ▶ Target validation: Network-based approaches to identify proteins whose modulation affects disease outcomes
- ▶ Combination therapy: Rational design of drug combinations based on complementary network effects

### Clinical Examples:

Aspirin's cardiovascular benefits arise from COX-1/COX-2 inhibition affecting multiple inflammatory pathways. Metformin activates AMPK but also affects mitochondrial complex I, explaining its pleiotropic effects. Imatinib targets BCR-ABL in chronic myeloid leukemia but also inhibits c-KIT and PDGFR, enabling use in gastrointestinal stromal tumors.

## • 5. Network Pharmacology

Network pharmacology represents a paradigm shift from the traditional "one drug, one target" approach to a systems-level understanding of drug action. This approach integrates data from genomics, proteomics, metabolomics, and clinical outcomes to design therapies that modulate disease networks rather than individual targets. It enables rational drug combination design, prediction of side effects, and identification of novel therapeutic opportunities.



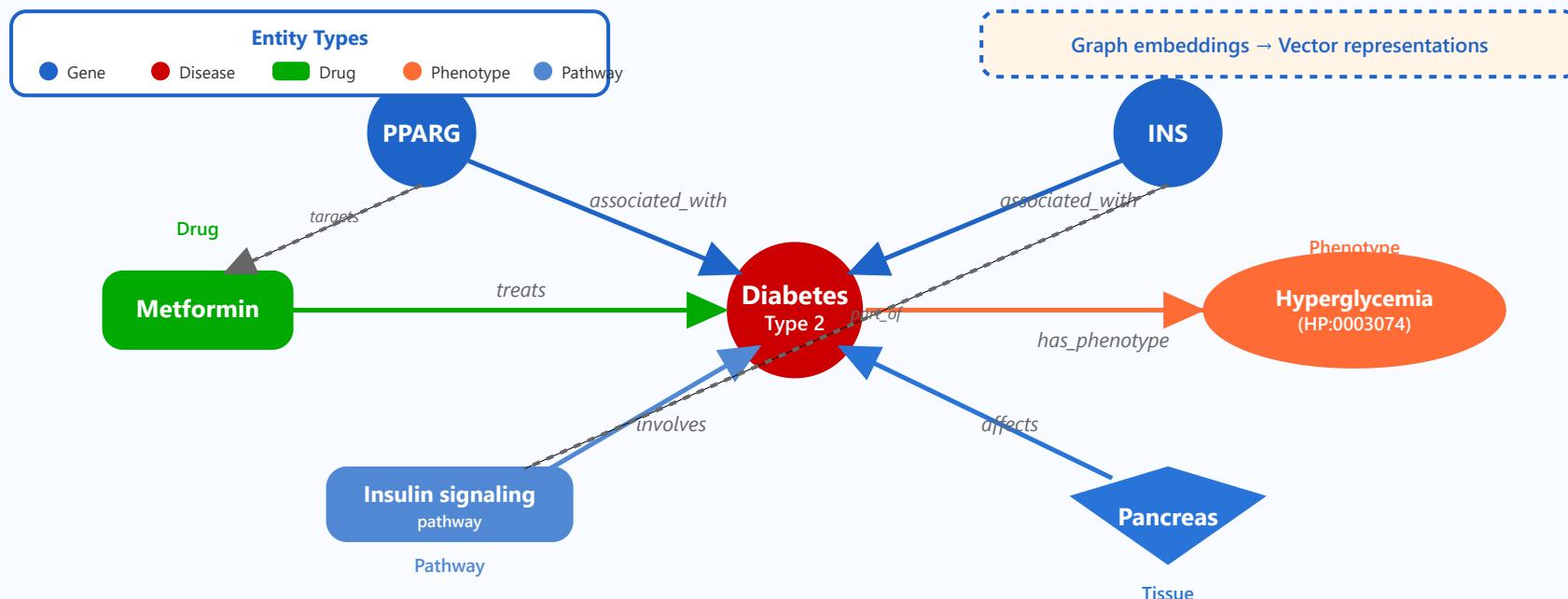
**Key Concepts:**

- ▶ Network-based drug discovery: Using network topology to identify optimal intervention points
- ▶ Combination therapy design: Rational selection of drug pairs that target complementary network components
- ▶ Drug repurposing: Identifying new therapeutic applications by analyzing drug-disease network proximity
- ▶ Personalized medicine: Customizing treatments based on patient-specific network perturbations
- ▶ Systems toxicology: Predicting adverse effects by analyzing drug impacts on cellular networks

### Clinical Examples:

Triple therapy for hypertension targets multiple pathways in the renin-angiotensin-aldosterone system. Cancer immunotherapy combinations like anti-PD-1 plus anti-CTLA-4 work synergistically by targeting different immune checkpoint pathways. Traditional Chinese Medicine formulations have been reinterpreted through network pharmacology, revealing multi-target mechanisms that explain their therapeutic effects.

# Biomedical Knowledge Graphs



## Key Components and Detailed Explanations

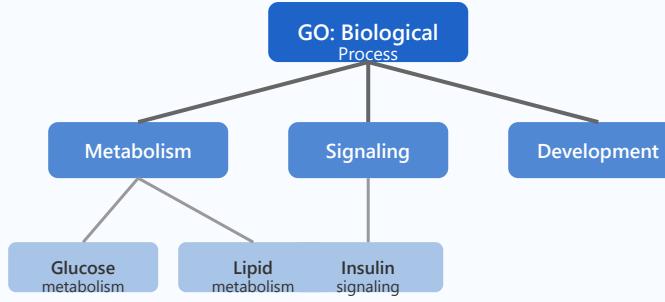
1

Biomedical Ontologies (생의학 온톨로지)

Biomedical ontologies are structured vocabulary systems that systematically define relationships between biological concepts. They enable consistent and standardized representation of biomedical knowledge.

#### Major Ontologies:

- ▶ **GO (Gene Ontology)**: Gene functions, biological processes, cellular components
- ▶ **HPO (Human Phenotype Ontology)**: Human phenotypes and disease characteristics
- ▶ **DO (Disease Ontology)**: Disease classification and relationships
- ▶ **ChEBI**: Chemical entities and drug information



Hierarchical Structure

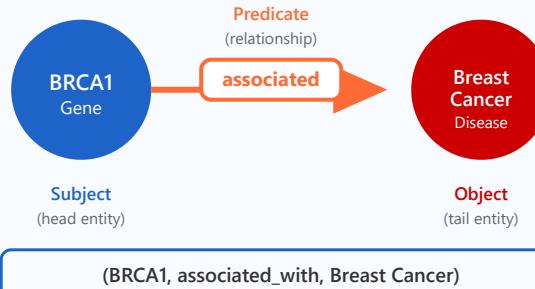
Connected by is-a relationships

## 2 Entity Relationships (개체 관계)

The core of knowledge graphs is explicitly representing relationships between various biomedical entities (genes, diseases, drugs, phenotypes, etc.). These relationships are utilized to discover new biological insights.

#### Major Relationship Types:

- ▶ **Gene-Disease**: BRCA1 gene mutation → Increased breast cancer risk
- ▶ **Drug-Target**: Imatinib → BCR-ABL kinase inhibition



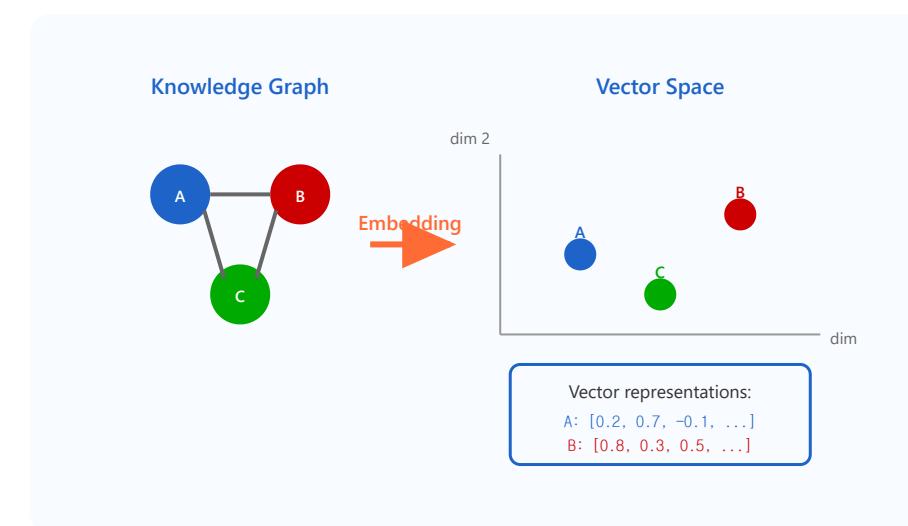
- ▶ **Disease-Phenotype:** Diabetes → Hyperglycemia
- ▶ **Gene-Pathway:** TP53 → p53 signaling pathway

### 3 Graph Embeddings (그래프 임베딩)

Graph embeddings transform nodes and edges of knowledge graphs into low-dimensional vector spaces. This enables machine learning models to effectively learn and utilize graph structure information.

#### Major Embedding Techniques:

- ▶ **TransE:**  $h + r \approx t$  (head + relation  $\approx$  tail)
- ▶ **Node2Vec:** Random walk-based node embeddings
- ▶ **GCN:** Graph Convolutional Networks
- ▶ **GAT:** Graph Attention Networks

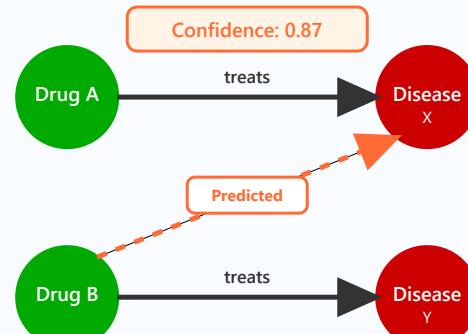


### 4 Link Prediction (링크 예측)

Link prediction is a technique for predicting yet undiscovered relationships between entities in knowledge graphs. It is utilized in drug development, disease mechanism elucidation, and gene function prediction.

#### Application Cases:

- ▶ **Drug Repurposing:** Discovering new indications for existing drugs
- ▶ **Target Identification:** Predicting drug target proteins
- ▶ **Disease Gene Discovery:** Finding disease-related genes
- ▶ **Adverse Effect Prediction:** Predicting drug side effects



Predicting that Drug B may also treat Disease X

## 5 Query Systems (쿼리 시스템)

Biological question answering systems allow researchers to ask complex biomedical questions in natural language and provide answers by retrieving relevant information from knowledge graphs. They utilize graph query languages like SPARQL.

#### Query Examples:

- ▶ "Find genes associated with Alzheimer's disease"
- ▶ "What drugs target EGFR?"
- ▶ "What are diabetes treatments and their side effects?"

#### Natural Language Query

"What drugs treat diabetes?"

#### SPARQL Query Translation

```
SELECT ?drug WHERE {  
?drug treats "Diabetes" }
```

#### Query Results

- ✓ Metformin - First-line treatment
- ✓ Insulin - Insulin replacement
- ✓ Sulfonylureas - Stimulate insulin secretion

- ▶ "Which pathways does TP53 gene participate in?"

# Case Studies in Multi-Omics and Systems Biology

## TCGA Pan-cancer

The Cancer Genome Atlas multi-omics integration

## METABRIC

Molecular taxonomy of breast cancer

## LINCS

Library of Integrated Network-based Cellular Signatures

## HuBMAP

Human BioMolecular Atlas Program

## Clinical Examples

Real-world clinical integration

1

## TCGA Pan-cancer Analysis

The Cancer Genome Atlas - Multi-Omics Integration

### ► Overview

The Cancer Genome Atlas (TCGA) represents one of the most comprehensive cancer genomics programs, analyzing over 33 different cancer types across more than 11,000 patients. The pan-cancer analysis integrates multiple omics

### Multi-Omics Data Integration

Genomics (WGS/WES)



layers to identify common molecular features and driver events across cancer types.

## ► Key Findings

- ✓ Identified 299 cancer driver genes across all tumor types
- ✓ Discovered shared molecular pathways across different cancers
- ✓ Characterized tumor microenvironment signatures
- ✓ Established molecular subtypes for precision medicine
- ✓ Created comprehensive mutation and expression landscapes

## Impact on Precision Medicine

TCGA data has enabled the development of tumor-agnostic therapies, where treatment is based on molecular features rather than tissue of origin, revolutionizing cancer treatment approaches.

Transcriptomics (RNA-seq)



Epigenomics  
(Methylation)



Proteomics (RPPA)



Integrated Analysis

## TCGA Statistics

**33**

Cancer Types

**11,000+**

Patients

**299**

Driver Genes

**2.5PB**

Data Generated

## ► Study Design

The Molecular Taxonomy of Breast Cancer International Consortium (METABRIC) analyzed nearly 2,000 breast cancer patients, integrating genomic and transcriptomic data to refine breast cancer classification and improve prognostic accuracy beyond traditional clinical markers.

## ► Major Discoveries

- ✓ Identified 10 integrative molecular subtypes (IntClust 1-10)
- ✓ Refined PAM50 intrinsic subtypes with genomic data
- ✓ Discovered copy number-driven subtypes
- ✓ Identified novel driver alterations in specific subtypes
- ✓ Improved prognostic stratification for treatment decisions

## Clinical Translation

METABRIC classifications are now used in clinical decision-making tools and have influenced the design of subtype-

## METABRIC Analysis Pipeline

Copy Number Arrays



Gene Expression Profiles



Clinical Annotations



Integrative Clustering



10 IntClust Subtypes

## Study Metrics

1,980

10

specific clinical trials, particularly for triple-negative and HER2-positive breast cancers.

Patients

IntClust Subtypes

5

PAM50 Subtypes

15y

Follow-up Period

3

## LINCS Program

Library of Integrated Network-based Cellular Signatures

### ► Program Goals

The Library of Integrated Network-based Cellular Signatures (LINCS) is an NIH program that catalogs how cells respond to various perturbations including genetic manipulations and chemical compounds. The goal is to understand cellular pathways and discover new therapeutic targets.

### ► Key Components

- ✓ L1000 gene expression signatures from perturbations
- ✓ Protein-protein interaction networks

### LINCS Workflow

Perturbation (Drug/Gene)



Cellular Response  
Measurement



Signature Generation



Database Integration

- ✓ Cell morphology profiling (Cell Painting)
- ✓ Transcription factor binding data
- ✓ Drug repurposing and combination predictions

## Therapeutic Prediction

### Applications in Drug Discovery

LINCS data enables systematic drug repurposing by identifying compounds with similar or opposite signatures to disease states, accelerating the discovery of new therapeutic applications for existing drugs.

### LINCS Database Scale

**1.3M+**

Perturbations

**25,000+**

Compounds

**70+**

Cell Lines

**978**

Landmark Genes

4

## HuBMAP

Human BioMolecular Atlas Program

### ► Program Mission

The Human BioMolecular Atlas Program (HuBMAP) aims to create a comprehensive, high-resolution molecular atlas of the human body at single-cell resolution. This ambitious project integrates spatial transcriptomics, proteomics,

### HuBMAP Data Integration

#### Tissue Collection

metabolomics, and imaging to map cellular architecture and function across organs.

### ► Technical Approaches

- ✓ Single-cell and spatial transcriptomics (MERFISH, seqFISH)
- ✓ Multi-modal imaging (IMC, CODEX, MIBI)
- ✓ Tissue mapping and 3D reconstruction
- ✓ Integration of molecular and anatomical data
- ✓ Open-access data portal and visualization tools

### Impact on Disease Understanding

By mapping healthy tissue architecture at unprecedented resolution, HuBMAP provides a reference atlas for understanding disease-related alterations in tissue organization, cellular composition, and molecular signaling.

Multi-Modal Profiling



Spatial Registration



3D Reconstruction



Interactive Atlas

### Program Scope

**80+**

Organ Types

**10+**

Assay Types

**1000s**

Tissue Samples

**Single-cell**

Resolution

## Clinical Integration Examples

## Precision Oncology Programs

Clinical implementation of multi-omics profiling has transformed cancer care through molecular tumor boards, where genomic, transcriptomic, and proteomic data inform treatment decisions. Programs like MSK-IMPACT and Foundation Medicine's comprehensive genomic profiling guide targeted therapy selection.

## Clinical Success Stories

- ✓ NTRK fusion detection leading to tumor-agnostic larotrectinib approval
- ✓ Liquid biopsy for minimal residual disease monitoring
- ✓ Pharmacogenomics guiding drug dosing (CYP2D6, TPMT)
- ✓ Multi-omic subtyping in acute myeloid leukemia (AML)
- ✓ Immune profiling for immunotherapy patient selection

## Future Directions

The integration of multi-omics data into electronic health records, real-time clinical decision support systems, and AI-

## Clinical Workflow Integration

Patient Sample Collection



Multi-Omics Profiling



Bioinformatics Analysis



Molecular Tumor Board



Treatment Decision

## Clinical Impact Metrics

30-40%

Actionable Findings

15-20%

Treatment Changes

driven treatment recommendations represents the next frontier in precision medicine implementation.

**10-15%**

Clinical Trial Matches

**2-4 weeks**

Turnaround Time

# Challenges in Multi-Modal Integration

## Missing Data

Incomplete measurements across modalities

## Batch Effects

Technical variation across platforms

## Scale Differences

Different measurement scales and distributions

## Interpretability

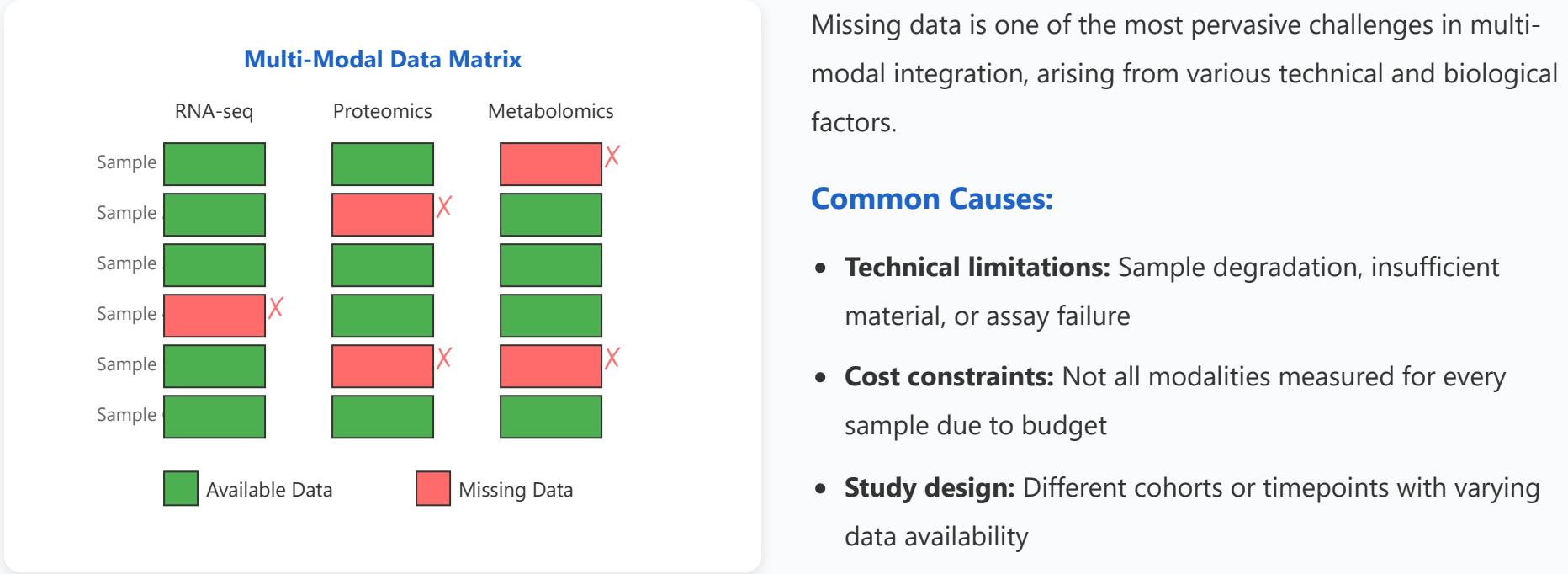
Understanding integrated models

## Validation

Reproducibility and generalization

1

## Missing Data



Missing data is one of the most pervasive challenges in multi-modal integration, arising from various technical and biological factors.

### Common Causes:

- **Technical limitations:** Sample degradation, insufficient material, or assay failure
- **Cost constraints:** Not all modalities measured for every sample due to budget
- **Study design:** Different cohorts or timepoints with varying data availability
- **Quality control:** Data filtered out due to quality metrics

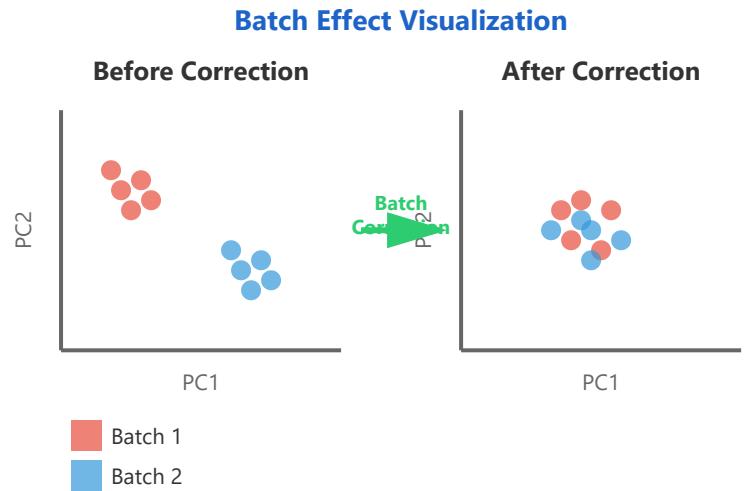
### Impact on Analysis:

- Reduced statistical power and sample size
- Biased results if missing data is not random (MNAR)
- Inability to apply certain integration methods requiring complete data
- Challenges in machine learning model training

**Solution Approaches:** Multiple imputation methods, pattern mixture models, methods robust to missing data (e.g., MOFA), or complete-case analysis with careful consideration of bias.

## 2

## Batch Effects



Batch effects represent systematic technical variation that can obscure true biological signals and lead to false discoveries in multi-modal integration.

### Sources of Batch Effects:

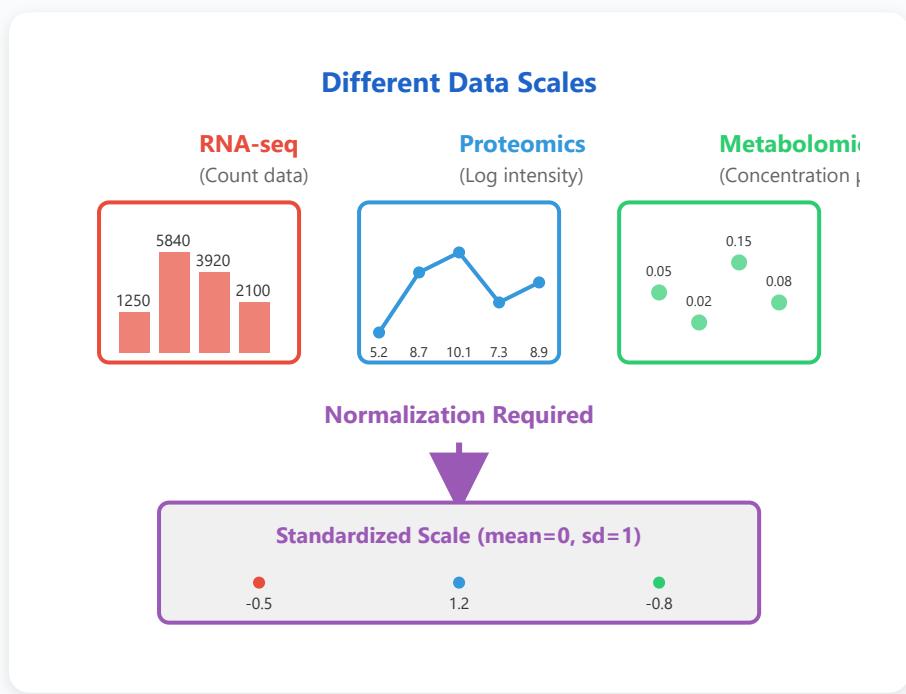
- **Laboratory differences:** Variations in protocols, equipment, and operators
- **Temporal effects:** Changes over time in reagents, instruments, or conditions
- **Platform differences:** Different measurement technologies or versions
- **Sample processing:** Storage conditions, extraction methods, handling time

### Consequences:

- Artificial clustering by batch rather than biology
- Inflated false positive rates in differential analysis
- Confounding with biological variables of interest
- Poor model generalization across studies

**Solution Approaches:** ComBat, limma's  
removeBatchEffect, Harmony, mutual nearest neighbors

### 3 Scale Differences



Different omics modalities measure distinct biological entities using various technologies, resulting in vastly different data scales and distributions that must be harmonized for integration.

#### Types of Scale Differences:

- Measurement units:** Counts (RNA-seq) vs. intensities (proteomics) vs. concentrations (metabolomics)
- Dynamic range:** Orders of magnitude difference in value ranges
- Distributions:** Negative binomial (RNA-seq), log-normal (proteomics), various (metabolomics)
- Sparsity:** Different proportions of zero or missing values

#### Integration Challenges:

- High-scale modalities dominating analysis without normalization

(MNN), or including batch as a covariate in statistical models. Proper experimental design with randomization is crucial.

- Invalid statistical assumptions when combining raw data
- Difficulty in defining meaningful distance metrics
- Feature weighting issues in machine learning models

**Solution Approaches:** Z-score normalization, quantile normalization, rank-based methods, variance stabilizing transformations (VST), or using methods that handle different scales internally (e.g., kernel-based approaches).

4

## Interpretability

### Model Interpretability Spectrum

**High Interpretability**

Linear Models

Factor Analysis

- ✓ Clear feature weights
- ✓ Biological insights
- ✓ Hypothesis generation

**Low Interpretability**

Deep Learning

Ensemble Methods

- X "Black box" nature
- X Hard to explain predictions

Often higher complexity = better performance but lower interpretability

Interpretability refers to the ability to understand and explain how integrated multi-modal models make predictions or identify patterns, which is crucial for generating biological insights and clinical trust.

### Why Interpretability Matters:

- **Biological discovery:** Understanding which features drive outcomes reveals mechanisms
- **Clinical application:** Healthcare decisions require explainable predictions
- **Model validation:** Detecting spurious correlations and biases
- **Regulatory requirements:** Medical applications often require interpretable models

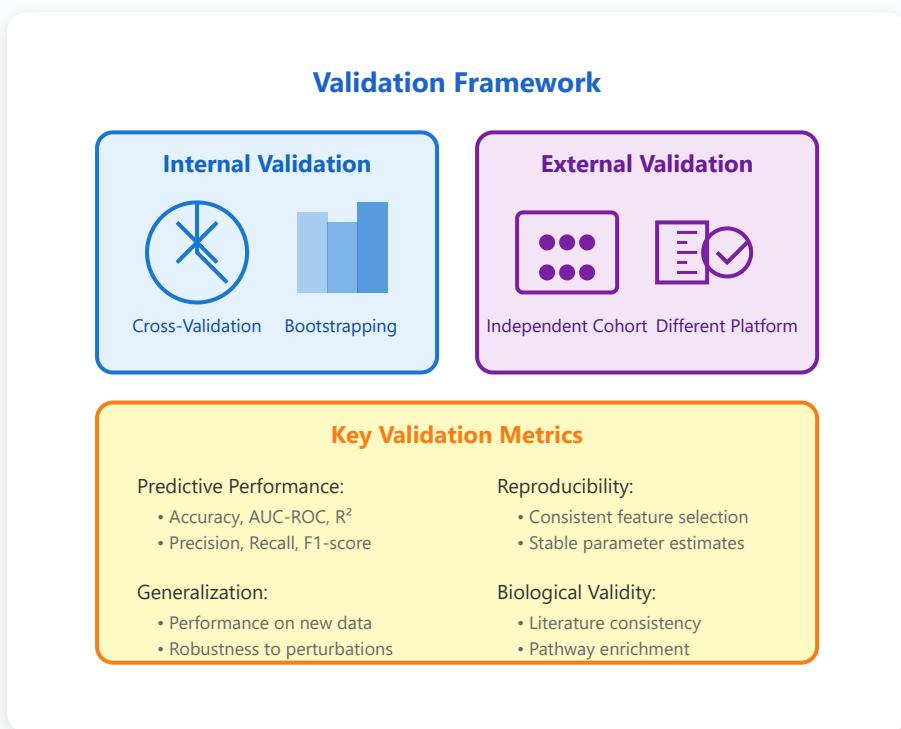
### Challenges in Multi-Modal Context:

- Complex interactions between modalities are hard to visualize
- High dimensionality obscures individual feature contributions
- Non-linear relationships complicate interpretation
- Trade-off between predictive accuracy and interpretability

**Solution Approaches:** SHAP values for feature importance, attention mechanisms in neural networks,

sparse models (LASSO, elastic net), factor analysis with loadings interpretation, or post-hoc explanation methods like LIME.

## 5 Validation



Validation ensures that multi-modal integration results are reliable, reproducible, and generalizable to new data, which is essential for translating findings into clinical applications or biological knowledge.

### Validation Challenges:

- Limited sample sizes:** Multi-modal datasets are often small, limiting validation power
- Overfitting risk:** High-dimensional data increases risk of spurious patterns
- Lack of standards:** No universal validation framework for multi-modal methods
- Cost constraints:** Independent validation cohorts are expensive to generate

### Best Practices:

- Use nested cross-validation for hyperparameter tuning and performance estimation
- Test on truly independent external cohorts when possible
- Validate biological findings through orthogonal methods or databases
- Assess stability through repeated subsampling or bootstrapping
- Report multiple performance metrics appropriate for the task

**Solution Approaches:** Rigorous cross-validation schemes, external validation cohorts, simulation studies with known ground truth, biological validation through experiments, and comprehensive reporting following established guidelines (e.g., TRIPOD for prediction models).

# Hands-on: MOFA (Multi-Omics Factor Analysis)

A Comprehensive Guide to Multi-Omics Integration

## Data Preparation

Formatting multi-omics datasets

## Model Training

Running MOFA analysis

## Factor Interpretation

Understanding learned factors

## Variance Decomposition

Attributing variance to factors

## Downstream Analysis

Using factors for prediction

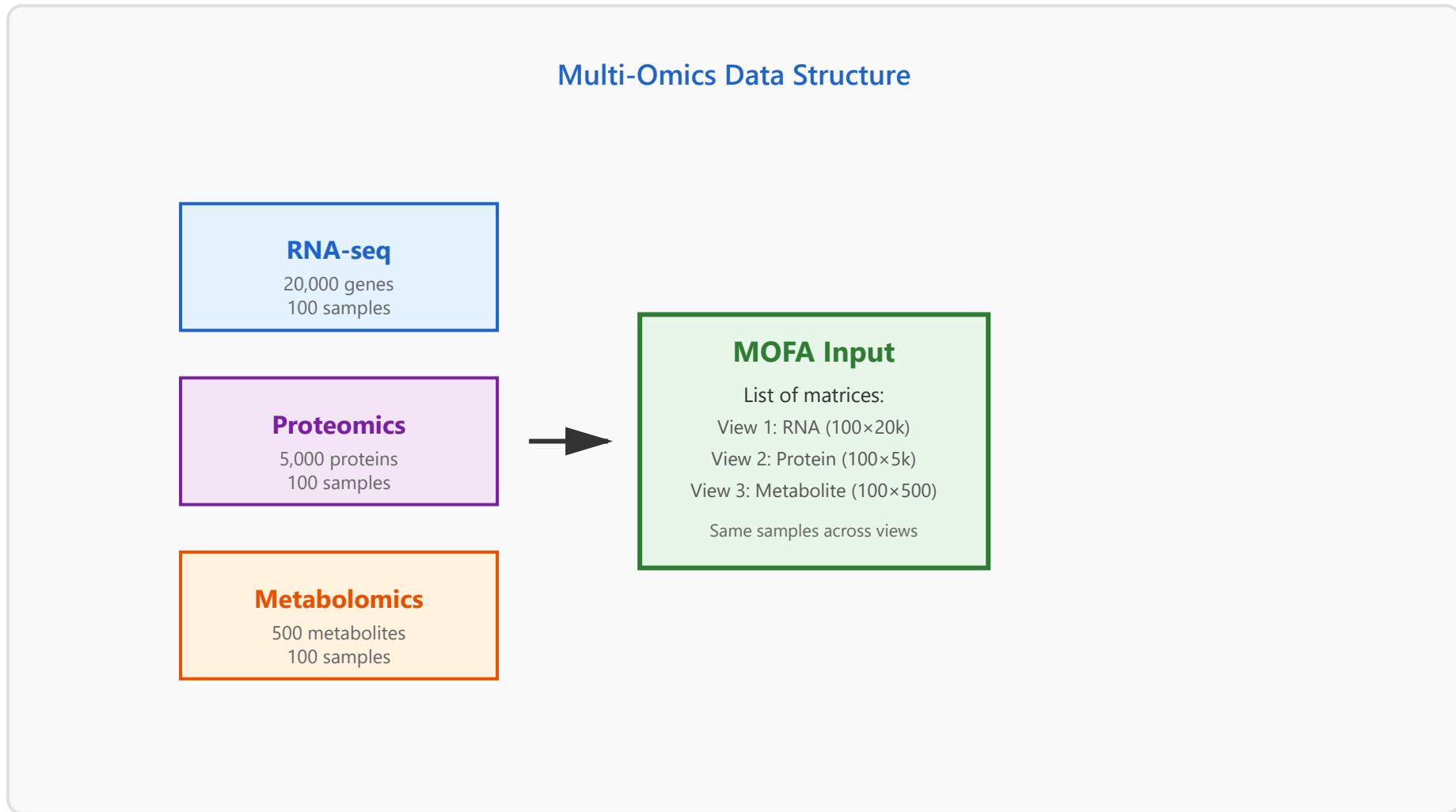
1

## Data Preparation

Data preparation is the foundational step in MOFA analysis. It involves organizing multiple omics datasets (e.g., genomics, transcriptomics, proteomics, metabolomics) into a structured format that MOFA can process. The key is to ensure that all datasets share the same samples while potentially having different features.

**Key Requirements:**

- **Sample alignment:** All omics layers must have measurements for the same set of samples
- **Data format:** Typically matrices where rows are samples and columns are features
- **Normalization:** Each omics layer should be appropriately normalized (log-transformation, standardization, etc.)
- **Missing data:** MOFA can handle missing values, but extreme missingness should be addressed
- **Feature selection:** Remove low-variance or non-informative features to improve computational efficiency



```
# Example R code for data preparation library(MOFA2) # Load your omics data rna_data <-  
read.csv("rna_seq.csv", row.names=1) protein_data <- read.csv("proteomics.csv", row.names=1)  
metabolite_data <- read.csv("metabolomics.csv", row.names=1) # Ensure samples match across datasets  
common_samples <- Reduce(intersect, list(rownames(rna_data), rownames(protein_data),  
rownames(metabolite_data))) # Create a list of matrices data_list <- list( "RNA" =  
as.matrix(rna_data[common_samples, ]), "Protein" = as.matrix(protein_data[common_samples, ]),  
"Metabolite" = as.matrix(metabolite_data[common_samples, ]) ) # Create MOFA object MOFAobject <-  
create_mofa(data_list)
```

### Key Points

- Always verify sample matching across all omics layers before analysis
- Consider filtering features with high missingness (>50%)
- Apply appropriate transformations (log for count data, scaling for continuous data)
- Document preprocessing steps for reproducibility

## 2 Model Training

Model training in MOFA involves learning latent factors that capture coordinated variation across multiple omics layers. MOFA uses a Bayesian framework with automatic relevance determination to identify the optimal number of factors and their relevance to each omics view.

### Training Process:

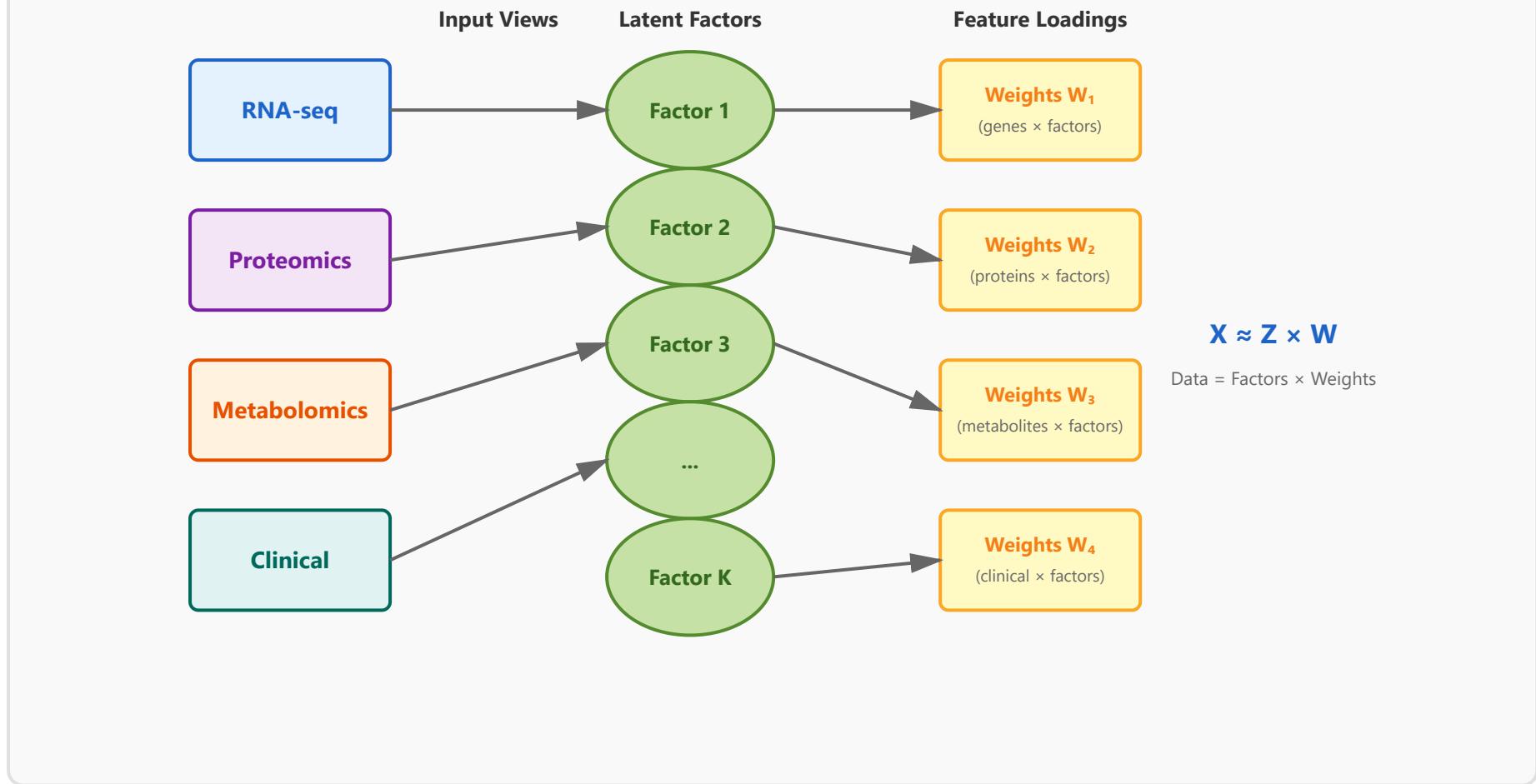
- **Factor initialization:** Start with more factors than expected (MOFA will prune irrelevant ones)

- **Iterative optimization:** Alternating updates of factor values and loadings
- **Convergence criteria:** Stop when the change in ELBO (Evidence Lower BOund) is small
- **Hyperparameter tuning:** Options for sparsity, number of factors, convergence thresholds

#### **Model Parameters:**

- **num\_factors:** Maximum number of factors (typically 10-25)
- **convergence\_mode:** "fast", "medium", or "slow" convergence
- **sparsity:** Degree of feature sparsity in loadings
- **seed:** Random seed for reproducibility

#### **MOFA Model Architecture**



```
# Configure and train MOFA model # Set model options
model_opts <- get_default_model_options(MOFAobject)
model_opts$num_factors <- 15 # Maximum number of factors
model_opts$sparsity <- TRUE # Enable sparse loadings # Set training options
train_opts <- get_default_training_options(MOFAobject)
train_opts$convergence_mode <- "medium"
train_opts$seed <- 42 # For reproducibility
train_opts$maxiter <- 1000 # Prepare MOFA object
MOFAobject <- prepare_mofa( object = MOFAobject, model_options =
model_opts, training_options = train_opts ) # Run MOFA training
MOFAobject <- run_mofa(MOFAobject, outfile = "model.hdf5") # Training typically takes 5-30 minutes depending on data size
```

### Key Points

- Start with more factors than expected (15-25) - MOFA will automatically prune unnecessary ones
- Convergence can take 5-30 minutes for typical datasets (100 samples, 3-5 views)
- Save the trained model to avoid re-training
- Monitor convergence by checking the ELBO plot

## 3 Factor Interpretation

After training, each latent factor represents a coordinated pattern of variation across omics layers. Interpreting these factors involves examining factor values (scores) for samples and feature weights (loadings) to understand biological meaning.

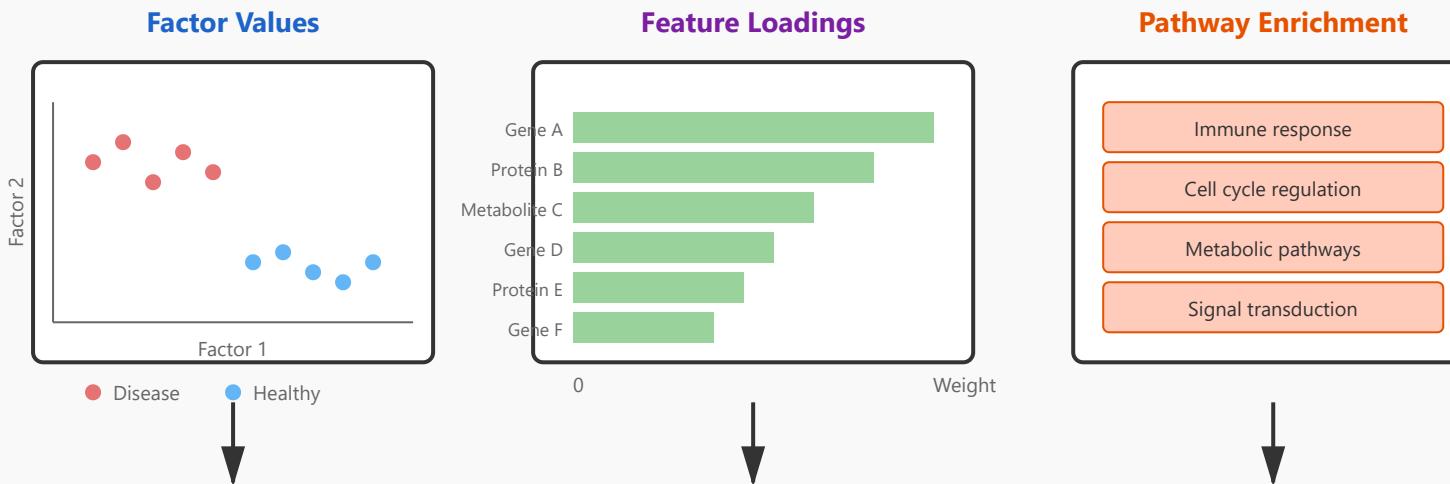
### Interpretation Approaches:

- **Factor values:** Examine how samples are distributed along each factor (e.g., by disease state, treatment group)
- **Feature weights:** Identify which genes/proteins/metabolites contribute most to each factor
- **Enrichment analysis:** Test for pathway/GO term enrichment in top weighted features
- **Correlation analysis:** Relate factors to clinical variables or phenotypes
- **View-specific patterns:** Determine which omics layers are most active in each factor

### Biological Insights:

- Factors may capture disease subtypes, treatment responses, or biological processes
- Multi-omics factors reveal coordinated regulation across molecular layers
- View-specific factors highlight layer-specific biological variation

# Factor Interpretation Workflow



## Biological Interpretation

### Factor 1: Disease Progression

- Separates disease vs. healthy samples
- High loadings on inflammatory genes and proteins
- Enriched for immune response pathways
- Correlates with disease severity scores
- Coordinated across transcriptome and proteome

```
# Extract and visualize factor values # Get factor values (Z matrix)
factors <- get_factors(MOFAobject,
factors = 1:5) # Plot factors colored by metadata
plot_factors(MOFAobject, factors = c(1,2), color_by =
"disease_status", shape_by = "treatment") # Get top weighted features
weights <- get_weights(MOFAobject,
views = "RNA", factors = 1, abs = TRUE)
top_genes <- head(sort(weights, decreasing=TRUE), 20) #

Visualize feature weights
plot_weights(MOFAobject, view = "RNA", factor = 1, nfeatures = 20) # Pathway
enrichment analysis
library(clusterProfiler)
gene_list <- names(top_genes)
enrichment <- enrichGO(gene =
```

```
gene_list, OrgDb = org.Hs.eg.db, ont = "BP", pAdjustMethod = "BH") # Correlate factors with clinical  
variables plot_factor_cor(MOFAobject, factors = 1:10, metadata = c("age", "BMI", "severity_score"))
```

### Key Points

- Always visualize factor values in the context of sample metadata (disease, treatment, etc.)
- Examine both positive and negative feature loadings for complete interpretation
- Use multiple sources of evidence (enrichment, literature, databases) to validate interpretations
- Not all factors need to have clear biological meaning - some capture technical variation

## 4 Variance Decomposition

Variance decomposition in MOFA quantifies how much variance in each omics layer is explained by each factor. This analysis reveals which factors are most important overall and which are view-specific or shared across views.

### Key Metrics:

- **R<sup>2</sup> per factor per view:** Proportion of variance explained in each omics layer by each factor
- **Total R<sup>2</sup>:** Overall variance explained by all factors in each view
- **Shared vs. specific factors:** Identify factors active across multiple views vs. view-specific
- **Factor importance ranking:** Order factors by total variance explained

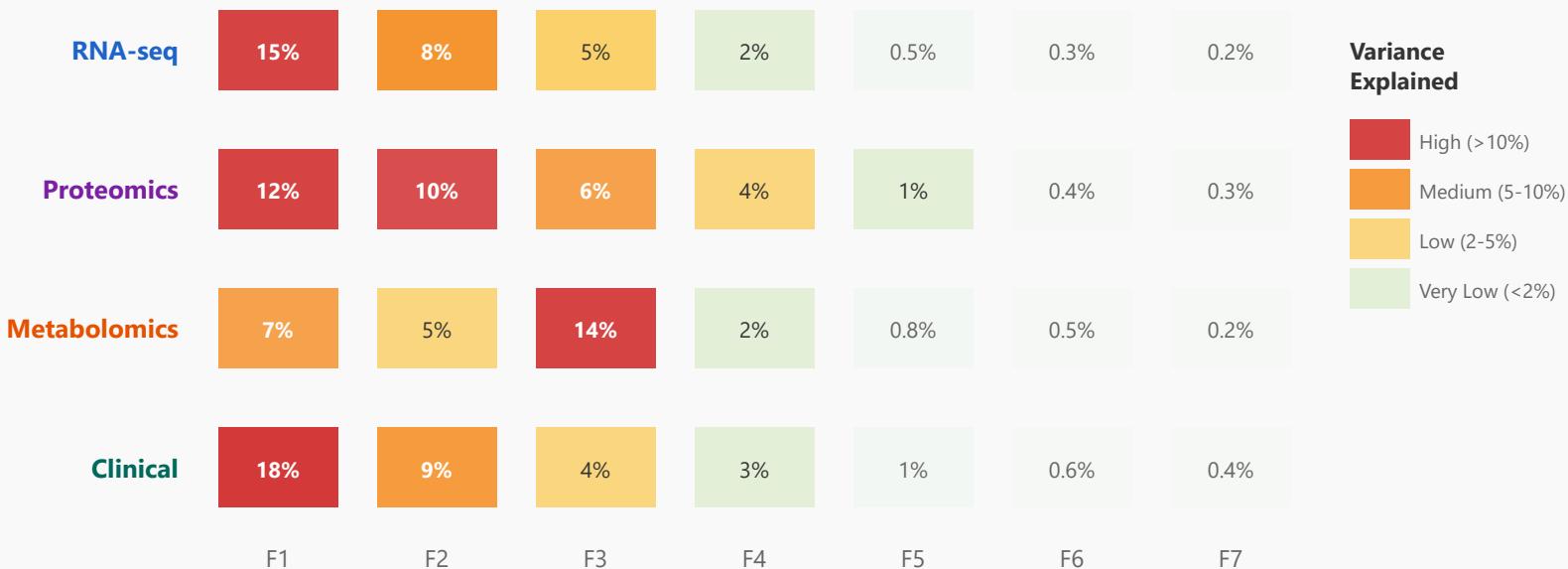
### Interpretation Guidelines:

- Factors explaining >5% variance in a view are typically considered important

- Shared factors (active in multiple views) represent coordinated biological processes
- View-specific factors capture layer-specific technical or biological variation
- Low total  $R^2$  may indicate noise, missing data, or need for more factors

## Variance Decomposition Heatmap

Variance Explained ( $R^2$ ) by Each Factor



### Key Observations

- is shared: high variance explained across all views (disease-related)
- is metabolomics-specific: captures metabolic variation
- Total variance explained: RNA (31%), Protein (34%), Metabolite (29%), Clinical (36%)

```
# Variance decomposition analysis # Calculate variance explained r2 <-
calculate_variance_explained(MOFAobject) # Plot variance explained by each factor in each view
plot_variance_explained(MOFAobject, x = "view", y = "factor") # Get total variance explained per view
total_r2 <- r2$r2_per_factor %>% group_by(view) %>% summarize(total_variance = sum(value))
print(total_r2) # view total_variance # RNA 0.31 # Proteomics 0.34 # Metabolomics 0.29 # Clinical 0.36 #
Identify factors with >5% variance in any view important_factors <- r2$r2_per_factor %>% filter(value >
0.05) %>% pull(factor) %>% unique() # Plot factor correlation structure plot_factor_cor(MOFAobject) #
Identify shared vs view-specific factors plot_factors_by_groups(MOFAobject, factors = 1:5, color_by =
"disease_status")
```

### Key Points

- Focus interpretation on factors explaining >5% variance in at least one view
- Shared factors (high R<sup>2</sup> across multiple views) are most biologically interesting
- View-specific factors may represent technical artifacts or layer-specific biology
- Low total R<sup>2</sup> (<20%) suggests factors may not fully capture data structure

## 5 Downstream Analysis

The latent factors learned by MOFA serve as powerful features for downstream machine learning and statistical analyses. These compressed representations capture coordinated multi-omics variation and can improve prediction, clustering, and classification tasks.

### Common Applications:

- **Predictive modeling:** Use factors as input features for predicting clinical outcomes (survival, response, disease progression)

- **Sample clustering:** Group samples based on factor profiles to identify subtypes
- **Dimensionality reduction:** Visualize samples in reduced factor space (PCA-like but multi-omics)
- **Biomarker discovery:** Identify factor-associated features as potential biomarkers
- **Integration with external data:** Relate factors to independent datasets or clinical variables

### **Advantages of Using MOFA Factors:**

- Reduced dimensionality (from thousands of features to 5-15 factors) improves model interpretability
- Factors integrate information across omics layers, capturing coordinated biology
- Less prone to overfitting compared to using raw high-dimensional data
- Biologically meaningful representations aid interpretation of predictions

### **Downstream Analysis Workflow**

Cox regression  
Kaplan-Meier curves

Random Forest  
SVM, Neural Networks

## MOFA Factors

Z matrix (100×10)

Compressed multi-omics

### Clustering

K-means, Hierarchical  
Subtype discovery

### Regression

Predict continuous  
outcomes (BMI, age)

AUC = 0.92

### Example: Disease Subtype Prediction

True Positive Rate



### Factor Importance

Factor 1	0.35
Factor 2	0.28
Factor 4	0.18
Factor 3	0.11
Factor 5	0.08

### Performance Metrics

Accuracy:  
Precision:  
Recall:  
F1-Score:

```
# Downstream analysis examples # 1. Extract factors as features factors_df <- get_factors(MOFAobject,  
as.data.frame = TRUE) metadata <- samples_metadata(MOFAobject) data_for_ml <- cbind(metadata,  
factors_df) # 2. Predictive modeling with Random Forest library(randomForest) library(caret) # Split  
data set.seed(42) trainIndex <- createDataPartition(data_for_ml$outcome, p=0.7, list=FALSE) train_data  
<- data_for_ml[trainIndex, ] test_data <- data_for_ml[-trainIndex, ] # Train model using MOFA factors  
rf_model <- randomForest(outcome ~ Factor1 + Factor2 + Factor3 + Factor4 + Factor5, data = train_data,  
ntree = 500, importance = TRUE) # Make predictions predictions <- predict(rf_model, test_data)
```

```
confusionMatrix(predictions, test_data$outcome) # 3. Survival analysis library(survival)
library(survminer) surv_data <- data_for_ml %>% select(Factor1, Factor2, Factor3, time, status) # Cox
proportional hazards model cox_model <- coxph(Surv(time, status) ~ Factor1 + Factor2 + Factor3, data =
surv_data) summary(cox_model) # Kaplan-Meier curves stratified by Factor 1 surv_data$Factor1_group <-
ifelse(surv_data$Factor1 > median(surv_data$Factor1), "High", "Low") fit <- survfit(Surv(time, status) ~
Factor1_group, data = surv_data) ggsurvplot(fit, pval = TRUE, risk.table = TRUE) # 4. Clustering for
subtype discovery factors_matrix <- as.matrix(factors_df) kmeans_result <- kmeans(factors_matrix,
centers = 3, nstart = 50) # Visualize clusters plot_factors(MOFAobject, factors = c(1,2), color_by =
kmeans_result$cluster, legend = TRUE) # 5. Correlation with clinical variables clinical_vars <- c("age",
"BMI", "tumor_size", "grade") cor_results <- cor(factors_df, metadata[, clinical_vars], use =
"pairwise.complete.obs") pheatmap(cor_results, cluster_rows = FALSE, cluster_cols = FALSE, color =
colorRampPalette(c("blue", "white", "red"))(100))
```

### Key Points

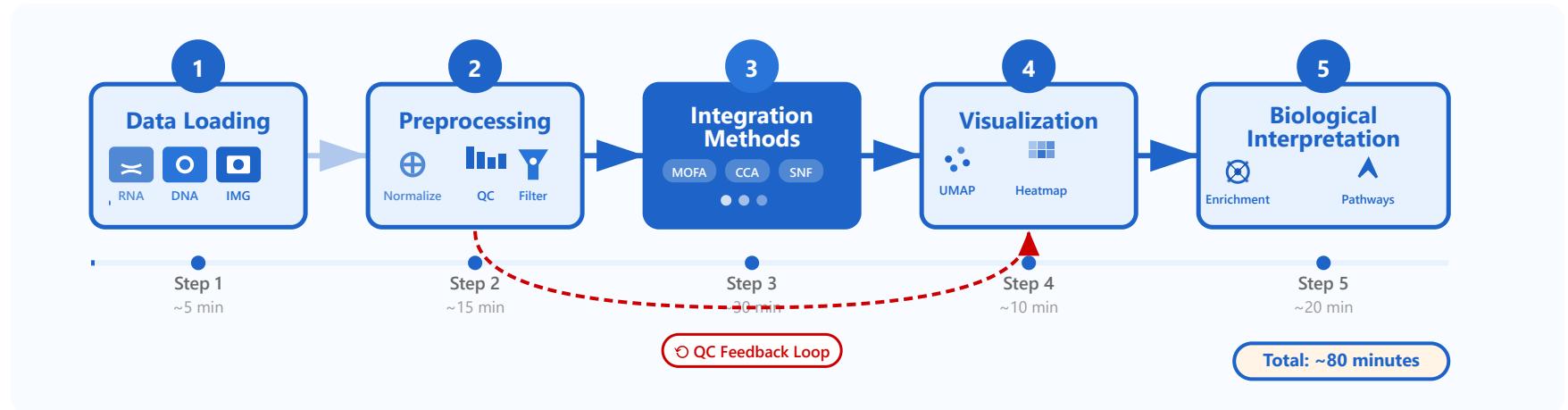
- MOFA factors provide interpretable, low-dimensional features that capture multi-omics variation
- Use cross-validation to assess predictive performance and avoid overfitting
- Factors often outperform single-omics data in prediction tasks due to information integration
- Combine MOFA with domain knowledge and validation in independent cohorts
- Factor-based models are more interpretable than black-box approaches on raw data

## Summary

MOFA provides a comprehensive framework for multi-omics integration, from data preparation through downstream analysis. By learning latent factors that capture coordinated variation across molecular layers, MOFA enables biological interpretation, variance

decomposition, and improved predictive modeling. The workflow presented here demonstrates the power of integrative approaches in understanding complex biological systems.

## Hands-on: Integration Workflow



### Data Loading

Reading multi-modal datasets

### Preprocessing

Normalization and quality control

### Integration Methods

Applying different integration approaches

### Visualization

UMAP, t-SNE, heatmaps

### Biological Interpretation

Functional enrichment analysis

# Thank you

Emerging methods in multi-modal integration  
Clinical impact and translational opportunities  
Research opportunities in systems medicine

**Introduction to Biomedical Data Science**