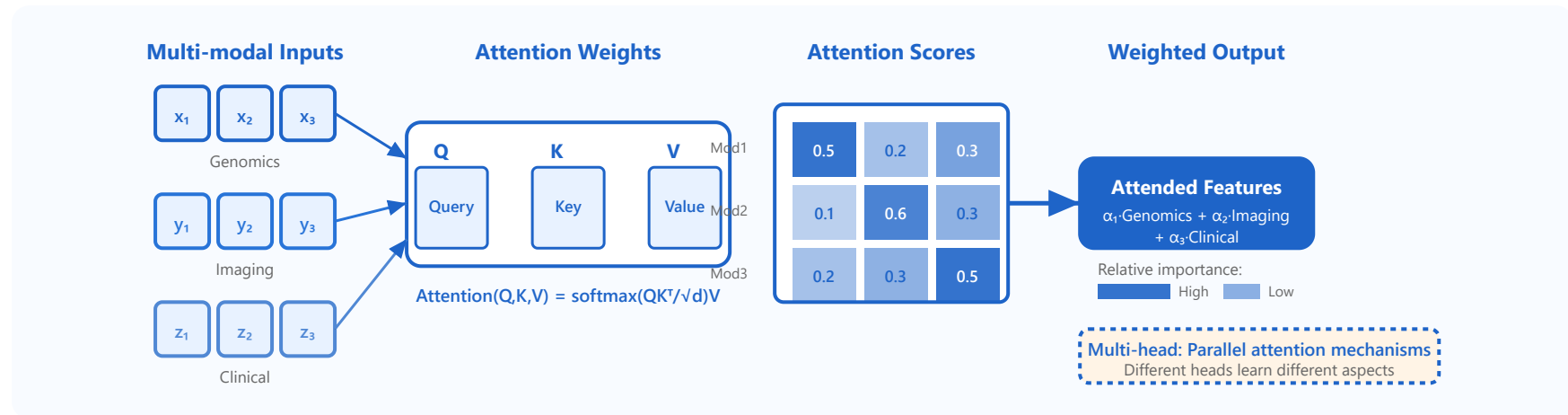


Attention-based Integration



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

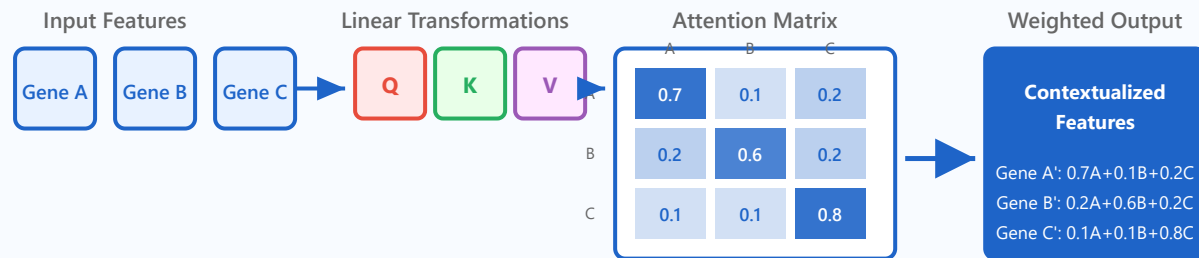
Overview: Self-attention fusion enables the model to learn dynamic importance weights for features within a single modality. Unlike fixed weighting schemes, self-attention adaptively determines which features are most relevant for the prediction task at hand.

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

Self-attention within Genomic Data



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

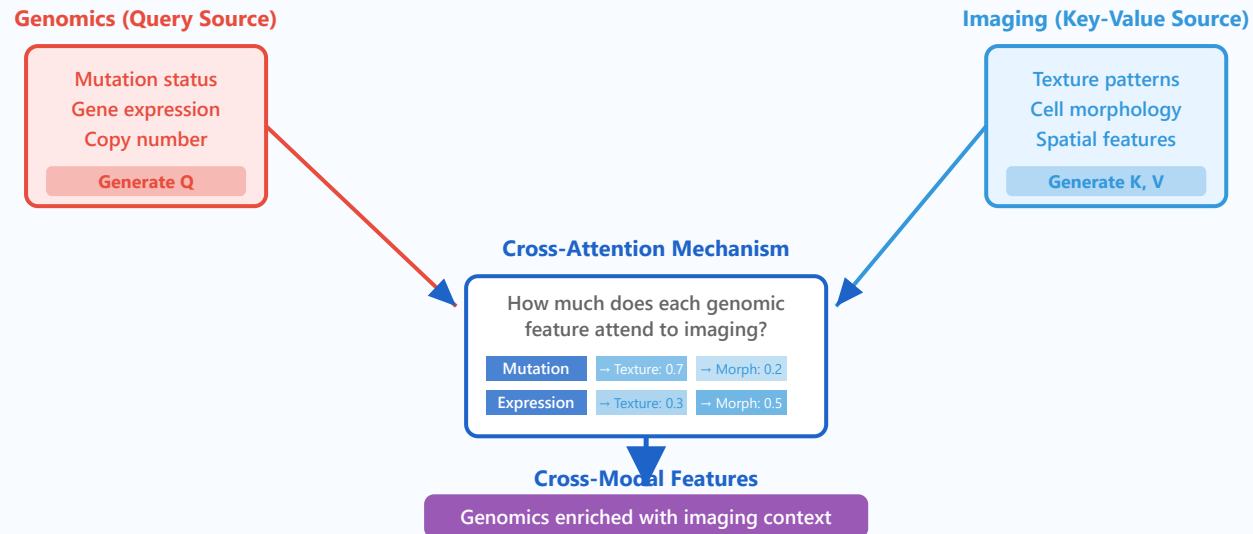
Overview: Cross-attention enables interaction between different data modalities by allowing one modality to attend to another. This mechanism is crucial for multi-modal learning where relationships between heterogeneous data sources must be discovered and leveraged.

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

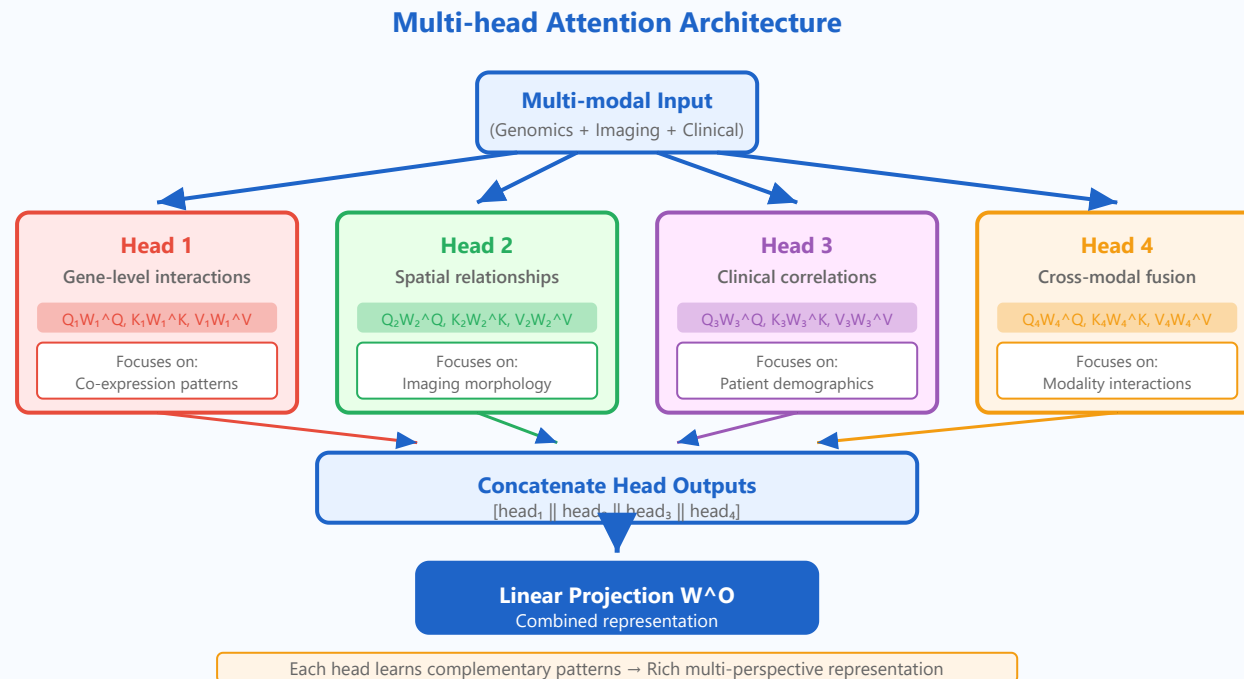
3. Multi-head Attention

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.

$$\text{MultiHead}(Q, K, V) = \text{Concat}(\text{head}_1, \dots, \text{head}_h) W^O$$

where $\text{head}_i = \text{Attention}(QW_i^Q, KW_i^K, VW_i^V)$



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

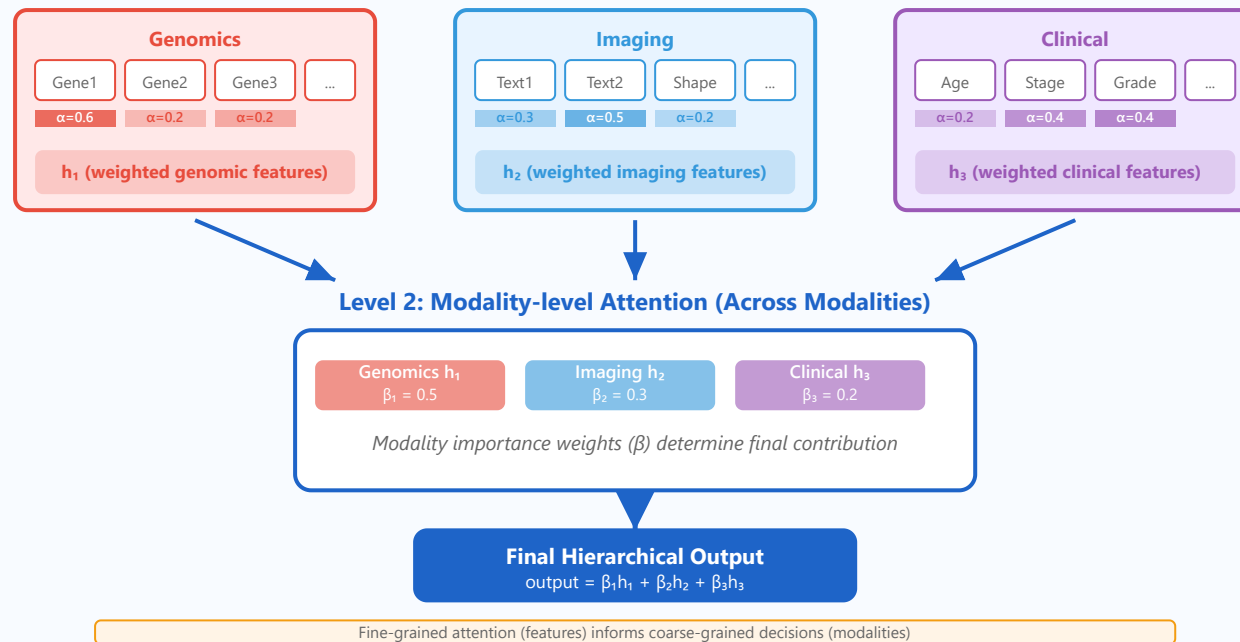
Overview: Hierarchical attention implements attention at multiple levels of granularity, typically combining feature-level attention with sample-level attention. This multi-scale approach mirrors the hierarchical structure of biological systems and clinical decision-making.

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

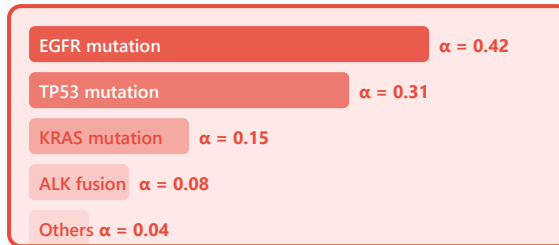
Overview: One of the most valuable aspects of attention mechanisms is their inherent interpretability. Attention weights provide direct insight into which features, samples, or modalities the model considers important for its predictions, making the decision-making process more transparent.

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.

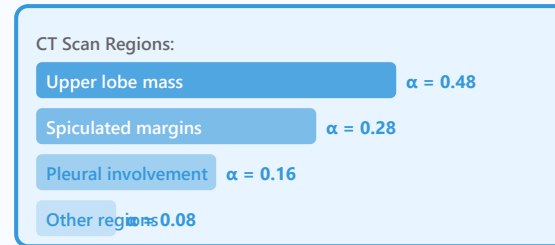
Attention-based Interpretability in Cancer Diagnosis

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