

# VLM 报告：特征解耦与解码

<https://gemini.google.com/share/80c8d4aee26b>

Responses below were generated with a creator's Gem according to their custom instructions. [Learn more](#) 5

March 2026 at 18:04

Published on 5 March 2026 at 18:10

1

5次+10000字+5公式 Custom Gem

这是一项极其宏大且极具前沿价值的研究任务。您提供的论文《Circuit Tracing in Vision-Language Models: Understanding the Internal Mechanisms of Multimodal Thinking》堪称人工智能“机制可解释性”(Mechanistic Interpretability) 领域的一座里程碑。它首次将目光从纯文本语言模型投向了更为复杂的视觉-语言多模态模型 (VLMs)，试图用类似神经生物学解剖的方式，逆向工程这些“黑盒”模型内部的思维回路。

按照您的要求，我将这份15000字的研究报告分为五个部分为您呈现。在**第一部分**中，我们将奠定整篇报告的理论基石，深入探讨VLM内部特征解耦的核心机制，并详细剖析该论文中最基础但也最重要的**前两个公式**（Transcoder的编码器与解码器公式）。

---

## 《多模态思维的内部机制：VLM回路追踪研究报告》(1/5)

—— 第一部分：打开黑盒与特征的升维解码

### 1. 引言：AI的“潜意识”与多义性迷局

近年来，诸如 CLIP、Flamingo、LLaVA 以及 Gemma-3 等视觉语言模型 (VLMs) 展现出了令人惊叹的能力。它们能够无缝地整合视觉感知与语言理解，不仅能回答关于图像的问题，甚至能跨模态进行复杂的多步推理。然而，伴随巨大成功而来的是一个令人不安的现实：这些模型依然是不透明的黑盒。我们知道它们“能做什么”，却根本不知道它们“在想什么”以及“为什么这么想”。在医疗影像诊断或自动驾驶等高风险领域，这种不透明性不仅限制了科学认知，更带来了难以预测的安全隐患。

在探索AI内部机制的过程中，研究人员遇到了一个被称为\*\*“多义性”

(Polysemanticity) \*\*的巨大阻碍。神经网络的参数量固然庞大，但它们需要学习和记忆的世界概念（例如颜色、形状、特定动物、抽象的数学符号等）要远多于模型内部的神经元数量。这就迫使模型采用一种“叠加”(Superposition) 的策略：**一个单一的神经元往往同时对多个毫无关联的概念做出响应。**

**具例说明：**假设模型中有一个编号为 #4096 的神经元。当你输入一张“海獭 (Sea Otter)”的照片时，它被高度激活；当你输入“蓝色 (Blue)”的文本时，它也被激活；甚至当你输入一个关于“微积分 (Calculus)”的公式时，它依然在放电。这种“一词多义”的神经元让我们根本无法追踪信息的真实流向。就像试图在一个交响乐团中，仅凭一个麦克风就听出第二小提琴手是否拉错了一个音符一样困难。

为了解开这个迷局，论文引入了\*\*Transcoder (转码器) \*\*技术，这也是我们今天要解析的最核心的两个公式的来源。

---

## 2. 核心突破：从多层感知机 (MLP) 到单义特征空间

在Transformer架构中，多层感知机 (MLP) 占据了模型约三分之二的参数，被认为是模型存储“知识”的记忆库。然而，MLP层包含着高度密集的非线性激活，使得传统的回路分析面临维度灾难。

以往的研究（如文本语言模型中的稀疏自编码器 SAE）旨在重建模型某一层的激活值。但Transcoder的设计更加巧妙：它直接替代了Transformer中的MLP子层，去学习并模拟整个MLP的输入-输出行为，同时保持计算的等效性。

通过训练Transcoder，研究人员成功地将密集、多义的残差流向量，映射到了一个极其宽广但高度稀疏的潜在空间中。在这个空间里，每一个特征 (Feature) 都变得\*\*“单义”(Monosemantic) \*\*——即一个特征只对应一个具体的、可解释的人类概念。

接下来，我们来看看实现这一“魔法”的两个核心公式。

---

## 3. 深度解析公式一：特征编码与激活 (Encoder)

Transcoder的第一步是将紧凑的、多义的输入向量展开。论文中给出了Transcoder编码器的公式：

$$z(x) = \text{ReLU}(W_{enc}x + b_{enc})$$

### 3.1 变量解构与物理意义

- $x \in \mathbb{R}^{d_{model}}$ : 这是MLP层的原始输入向量。它是一个包含着密集信息的低维向量（在Gemma-3-4B中， $d_{model} = 2560$ ）。此时的  $x$  就像是一个被极度压缩的“概念线团”。
- $W_{enc} \in \mathbb{R}^{d_{feat} \times d_{model}}$ : 这是编码器的权重矩阵。请注意它的维度， $d_{feat} \gg d_{model}$ （特征维度远大于输入维度）。这相当于提供了一个极其广阔的“高维画布”。
- $b_{enc} \in \mathbb{R}^{d_{feat}}$ : 编码器的偏置项，用于调节各个特征的激活阈值。
- $z(x)$ : 这是最终得到的**潜在特征激活值**（Latent Feature Activations）。
- $ReLU$ : 线性整流函数（Rectified Linear Unit），一种非线性激活函数。

### 3.2 机制洞察与具例说明

这个公式的深刻之处在于\*\*“高维展开”与“稀疏截断”的结合\*\*。

设想你正在要求VLM处理一张“海獭在水面上挥手”的图片。

当视觉Token的信号传递到MLP层时，输入向量  $x$  处于一种混沌状态。它同时包含了“水”、“哺乳动物”、“棕色”、“反光”等无数概念的混合。

1. **高维展开 ( $W_{enc}x$ )**: 权重矩阵将这个2560维的混合向量，强行投影到一个数万乃至数十万维的庞大空间中。在这个高维空间里，原本挤在一起的概念有了足够的“房间”被分离开来。这在数学上类似于支持向量机（SVM）中的核技巧（Kernel Trick）——在低维空间线性不可分的数据，投射到高维空间后往往就变得线性可分了。
2. **稀疏截断 ( $ReLU$  与  $Top-K$ )**: 并非所有的房间都需要住人。为了确保每个特征的纯粹性，我们只保留那些真正强烈相关的信号。在实际操作中，研究团队不仅使用了ReLU，还采用了  $TopK(z(x), k)$  策略（保留前  $k$  个最大的激活值，实验中  $k = 48$ ），强行将其余特征置为零。

**结果:** 经过公式一的洗礼，原本混沌的向量被分解为极少数明确发光的“灯泡”。在这个广阔的高维稀疏空间  $z(x)$  中，可能只有第 #10204 号特征亮了，而通过后续分析我们会发现，这个特征完美且唯一地对应着“太平洋海岸线（Pacific coastline）”或“海獭（Sea otters）”这一特定语义。我们终于从黑盒中提取出了人类可以理解的\*\*“单义特征”\*\*！

---

## 4. 深度解析公式二：特征的重构与输出（Decoder）

将信息解析为单义特征还不够，模型必须将这些概念重新“打包”，以便传递给下一层继续计算。这就引出了Transcoder的解码器公式：

$$TC(x) = W_{dec}z(x) + b_{dec}$$

#### 4.1 变量解构与物理意义

- $z(x)$ : 上一阶段得到的稀疏单义特征向量。
- $W_{dec} \in \mathbb{R}^{d_{model} \times d_{feat}}$ : 解码器的权重矩阵。它的作用是将高维的稀疏特征重新压缩回原始模型维度（例如2560维）。
- $b_{dec} \in \mathbb{R}^{d_{model}}$ : 解码器的偏置项。
- $TC(x)$ : Transcoder的最终输出。它的目标是无限逼近原始MLP层的输出（即  $TC(x) \approx MLP(x)$ ）。

#### 4.2 机制洞察与具例说明

公式二本质上是一个\*\*\*“字典查询与线性组合”\*\*\*的过程。

我们可以把  $W_{dec}$  的每一列看作是字典中的一页（一个特征的基向量）。因为  $z(x)$  是极度稀疏的（最多只有48个非零值），所以公式二的计算实际上是在巨大的字典中，挑出那48页有用的内容，然后按照它们激活值的强度（权重）将它们相加。

**具例说明：** 延续海獭的例子。经过公式一，我们得到了几个被激活的单义特征 ( $z(x)$ )，比如：

- 特征A（权重0.8）：代表“水生动物”。
- 特征B（权重0.6）：代表“毛茸茸”。
- 特征C（权重0.9）：代表“特定的地理位置（太平洋）”。

此时，解码器矩阵  $W_{dec}$  里面存有这些特征对应的“机器语言”代码。公式二  $TC(x)$  的工作，就是用  $0.8 \times [\text{水生动物的机器代码}] + 0.6 \times [\text{毛茸茸的机器代码}] + 0.9 \times [\text{太平洋的机器代码}] + \text{偏置项}$ 。

最终，这些概念又被无缝缝合回了一个2560维的残差流向量  $TC(x)$  中，继续流向Transformer的下一层。

通过这两个公式，研究者实现了极为壮观的壮举：**在不干扰模型正常前向传播的前提下，在每一层内部搭起了一个“显微镜”**。原本不可名状的矩阵乘法，被优雅地拆解为了“概念分离（Encoder）”与“概念重组（Decoder）”的过程。

---

## 5. 第一阶段总结与前瞻

通过上述两个公式（编码器与解码器），这篇论文解决了机制可解释性走向多模态领域的第一只拦路虎——特征解耦。跨越文本与图像两种截然不同统计特性的模态数据，终于在一个统一的、单义的稀疏潜在空间中找到了可量化的表征形式。

但解耦仅仅是第一步。提取出单义特征后，我们如何衡量这面“显微镜”的清晰度？更重要的是，当这些孤立的特征在神经网络的几十层之间互相传递时，它们是如何产生因果联系，进而形成逻辑推理的？为什么视觉上的“火星”图片，会在网络深层自动激活文本概念上的“航天飞机”？

在下一部分的研究报告中，我们将继续深入，重点解析论文中的公式三（FVU：未解释方差比例）以及公式四（Attribution：归因权重），探讨模型如何确保这种特征拆解的忠实度，并揭开构建“多模态因果归因图（Attribution Graphs）”的数学奥秘。

---

### 附录：外部补充材料与参考文献（持续更新）

1. Anthropic (2023). *Towards Monosemanticity: Decomposing Language Models With Dictionary Learning*. (解释了稀疏自编码器如何在单模态中分离出单义特征，为本论文的Transcoder提供了思想基础)。
2. OpenReview (2024). *Transcoders Find Interpretable LLM Feature Circuits*. (这篇 NeurIPS 会议论文详细讨论了 Transcoder 相较于传统 SAE 的优势，即直接逼近 MLP 输出而非输入，从而消除非线性带来的回路追踪障碍)。 URL: <https://arxiv.org/abs/2406.11944>
3. *A Technical Primer on Mechanistic Interpretability*. (综述了特征提取、回路发现以及普遍性这三大机制可解释性支柱)。

---

是否需要我继续为您生成第二部分（聚焦模型保真度评估与因果归因图的构建）？

### 《多模态思维的内部机制：VLM回路追踪研究报告》(2/5)

#### —— 第二部分：重构保真度评估与因果归因图的构建

在第一部分中，我们探讨了如何利用 Transcoder（转码器）的编码器和解码器公式，将 VLM 内部混沌、多义的隐藏状态，强行展开为一个广阔且稀疏的单义特征空

间。这就像是我们成功地将交响乐团中混合的声波，拆解成了单独的乐器音轨。

然而，科学研究不能仅仅停留在“看起来有效”的直觉上。我们面临着两个极其严峻的拷问：

- 保真度拷问：**在强制要求“单义”和“稀疏”（只保留 Top-48 激活值）的过程中，Transcoder 是否丢失了原始 MLP 层中关键的“暗物质（Dark Matter）”或非线性信息？
- 因果性拷问：**即使我们提取出了独立的“海獭”或“火星”特征，但孤立的特征并不是“思维”。在拥有数十层的神经网络中，第 15 层的某个视觉特征，是如何跨越残差流，精准触发第 20 层的某个语义特征的？

为了回答这两个问题，这篇论文引入了极其精妙的保真度评估与因果追踪机制。这就引出了我们今天要深入解构的**第三个和第四个核心公式**。

---

## 1. 深度解析公式三：未解释方差比例（FVU）

在机制可解释性（Mechanistic Interpretability）领域，我们必须量化“显微镜”的失真程度。论文使用了\*\*未解释方差比例（Fraction of Variance Unexplained, 简称 FVU）\*\*来评估 Transcoder 重构原始 MLP 行为的质量：

$$FVU = \frac{\frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2}{\frac{1}{n} \sum_{i=1}^n (y_i - \bar{y})^2} = \frac{MSE}{Var(y)}$$

### 1.1 变量解构与物理意义

- $y_i$ ：原始 VLM 模型中，多层感知机（MLP）在真实前向传播时的实际输出向量。这代表了“真实的基准（Ground Truth）”。
- $\hat{y}_i$ ：Transcoder 仅使用少量极其稀疏的单义特征，重新组合计算出的近似输出向量（即  $TC(x)$ ）。
- $\bar{y}$ ：所有原始输出样本的平均值。
- 分子 (MSE)：**均方误差（Mean Squared Error）。它衡量了 Transcoder 重构出来的信息，与原始信息之间到底差了多少。
- 分母 (Var(y))：**原始数据的总方差。它代表了 MLP 原始输出本身包含的信息波动总量。

### 1.2 机制洞察与具例说明

FVU 的直观含义非常粗暴且有效：“在原始模型产生的所有变化中，有多少百分比是你这个 Transcoder 无法解释的？” FVU 的值越接近 0，说明 Transcoder 越完美地复刻了原始模型的计算逻辑；FVU 越高，说明强行引入的稀疏性导致了严重的信息丢失。

**具例说明：** 假设原始模型的 MLP 层输出是一个包含极高细节的 4K 分辨率图像（代表  $y$  的总方差）。而我们的 Transcoder 为了追求可解释性，只用了 48 种基础颜色（特征）来重新绘制这幅画（代表  $\hat{y}$ ）。如果重绘出来的画，在结构和色彩上与原画高度一致（MSE 很小），那么 FVU 可能只有 0.05（即 95% 的信息都被这 48 种特征解释了）。反之，如果丢失了关键的阴影和渐变，FVU 就会飙升。

**令人振奋的实证发现：** 论文在针对 Gemma-3-4B-it 模型的实验中揭示了一个极具启发性的现象：当仅使用纯文本数据训练 Transcoder 时，模型中间层（如第 15 层）的 FVU 居高不下；但一旦引入了多模态数据（文本+图像）进行联合训练，FVU 出现了显著的下降。

这在理论上证明了极其深刻的一点：VLM 中间层的表征本质上是多模态交织的。试图仅用纯文本的逻辑去解释这些中间层，必然会遗漏大量的方差；只有同时用视觉和语义特征进行约束，才能真正解释模型的内部表征。

---

## 2. 从特征到网络：构建因果归因图 (Attribution Graphs)

即便 FVU 极低，证明了特征的保真度，但一堆零散的高质量零件并不能直接组成一台运转的发动机。思维是一个连贯的过程。

为了弄清楚特征是如何相互影响的，论文引入了\*\*归因图 (Attribution Graphs)\*\* 的概念。在给定特定提示词 (Prompt) 和图像的情况下，由于我们将所有的非线性操作（如 LayerNorm、注意力机制的 Softmax、ReLU）都“冻结”在了该输入下的固定值，整个神经网络在局部变成了一个纯线性的函数。

这使得我们可以完美地计算出一条“因果链”。这就是论文中的公式四（特征归因公式）：

$$A_{s \rightarrow t} = a_s w_{s \rightarrow t}$$

以及它背后隐藏的第五个核心公式（虚拟权重展开式）：

$$w_{s \rightarrow t} = f_{dec}^{(s)\top} J_{(s) \rightarrow (t)}^r f_{enc}^{(t)}$$

## 2.1 变量解构与“信使”的传递

- $A_{s \rightarrow t}$ : 这是从浅层网络（如第 5 层）的源特征  $s$ ，对深层网络（如第 15 层）的目标特征  $t$  的**直接归因值 (Attribution)**。它代表了  $s$  的存在，在多大程度上直接导致了  $t$  的激活。
- $a_s$ : 源特征  $s$  在当前输入下的实际激活强度。
- $w_{s \rightarrow t}$ : 这是一个**虚拟权重 (Virtual Weight)**。它是将相隔多层的两个特征强行联系起来的数学桥梁。
- $f_{dec}^{(s)\top}$ : 源特征  $s$  的解码器向量。它定义了特征  $s$  是如何将自己的信息“写入”神经网络的主动脉（残差流）中的。
- $J_{(s) \rightarrow (t)}^r$ : 这是一个雅可比矩阵 (Jacobian Matrix)。它囊括了从  $s$  层到  $t$  层之间所有的注意力机制和信息路由操作。你可以把它理解为信息在神经网络海洋中漂流的“洋流图”。
- $f_{enc}^{(t)}$ : 目标特征  $t$  的编码器向量。它定义了特征  $t$  是如何从残差流中“读取”信息的。

## 2.2 机制洞察与跨模态联想的具例

公式四和公式五的结合，堪称这篇论文中最暴力的美学体现。它将原本需要穷举消融 (Ablation) 才能验证的因果关系，转化为了一次优雅的线性代数点乘。

### 具例说明：火星与航天飞机的跨模态联想

论文中展示了一个令人拍案叫绝的发现。研究人员输入了一张纯粹的“火星”照片，并配上提示词“This is the planet \_”。

1. **源特征 (s) 激活**: 在较浅的网络层中，图像嵌入 (Image Embedding) 激活了一个纯视觉的特征  $s$ ——“这是一个红色的星球（视觉上的行星特征）”。此时， $a_s$  的值很高。
2. **写入与路由 ( $f_{dec}^{(s)\top} J$ )**: 这个“红色星球”的特征将其信号写入残差流。通过雅可比矩阵  $J$  代表的注意力机制，这个信号被精准路由到了序列的后续位置。
3. **深层特征读取 ( $f_{enc}^{(t)}$ )**: 在网络的中深层，存在一个特征  $t$ 。通过特征可视化，研究人员发现特征  $t$  竟然对应着\*\*“航天飞机 (Space Shuttle)”\*\*的视觉概念。
4. **完成归因**: 虽然输入的图片中根本没有航天飞机，也没有任何文字提到航天飞机，但公式四计算出的归因值  $A_{s \rightarrow t}$  异常强烈。

**深刻结论：**这证明了 VLM 在内部建立了一张**脱离了文本指令的纯视觉联想网络**。看到火星，模型的视觉潜意识回路就会自动唤醒“航天飞机”的特征。这与人类“触景生情”的联想机制惊人地相似。

---

### 3. 第二阶段总结与前瞻

借助 FVU 公式（保证了解码的高保真）与 Attribution 归因公式（刻画了特征间的因果作用力），我们终于拥有了绘制“AI 脑图”的完整工具链。这套方法论不仅在理论上自洽，更在实践中首次揭示了多模态模型内部“视觉与语义双线并行、最终高层收敛”的层级结构。

但是，绘制出复杂的归因网络图仅仅是解读模型的第一步。在实际应用中，这种网络往往包含数以千计的节点，充满了噪音和冗余的“微弱连接”。

在接下来的**第三部分**中，我们将离开纯粹的数学公式，走向更具实验性和挑战性的实证环节：**我们该如何利用人类专家的先验知识与注意力分析机制，去修剪这棵庞大的因果树？面对那些无法用文本轻易描述的“纯视觉特征”，研究人员又开发了怎样的破解之法？**

您准备好和我一起进入 VLM 视觉注意力热力图（Attention Maps）的解析环节了吗？（回复“继续”或提出您对这两个公式的具体疑问，我们立即开启下一部分的研究）

3

5次+10000字+5公式 Custom Gem

#### 《多模态思维的内部机制：VLM回路追踪研究报告》(3/5)

##### —— 第三部分：视觉潜意识的破译与注意力热力图的数学重构

欢迎回到我们的研究之旅！在前面的两部分中，我们通过解析前四个核心公式（编码器、解码器、FVU 与特征归因公式），成功地在 VLM 内部搭建了一架“高维显微镜”，并绘制出了特征之间相互作用的因果归因图（Attribution Graphs）。

但是，当我们真正拿着这幅庞大的归因图去审视 VLM（例如 Gemma-3-4B）时，一个极其棘手的问题出现了：**文本特征是会“说话”的，但视觉特征是“失语”的。**

当模型激活了一个文本特征时，我们可以通过查看哪些文本片段（如“太平洋”、“水生动物”）以高频激活它，从而轻易给它打上标签。然而，Gemma-3 的视觉编码器（SigLIP）会将一张  $896 \times 896$  的图像，压缩、池化为 256 个“软图像 Token（Soft

Image Tokens)”输入给语言模型。这 256 个 Token 是一堆毫无人类可读语义的密集向量。当归因图告诉我们“第 14 个视觉 Token 强烈触发了后续的推理”，我们依然一头雾水——这第 14 个 Token 到底“看”到了图像的哪里？是海獭的爪子，还是火星的陨石坑？

为了打破这种跨模态的语义壁垒，研究团队巧妙地将自然语言处理（NLP）领域的注意力展开技术（Attention Rollout）引入了视觉塔（Vision Tower）。这就引出了本报告的**第五个核心公式：注意力展开与聚合公式**。

## 1. 深度解析公式五：视觉注意力聚合（Attention Rollout）

为了将抽象的视觉特征重新映射回人类可理解的二维物理空间（即原始图像的像素区域），论文使用了一种多层注意力矩阵的连乘与归一化方法。其核心公式表达如下：

$$R = \tilde{A}^{(L-K+1)} \tilde{A}^{(L-K+2)} \dots \tilde{A}^{(L)}$$

### 1.1 变量解构与计算步骤

在揭示这个公式的物理意义之前，我们需要严谨地解构构建矩阵  $\tilde{A}^{(l)}$  的前置步骤：

- **挑选最专注的“眼睛”（低熵过滤）**：视觉编码器在第  $l$  层有多个注意力头（Attention Heads）。有些头是“散视”的（全局关注，高熵），有些头是“聚焦”的（死死盯着某个局部，低熵）。研究者设定了一个比例  $q$ ，**仅挑选出平均熵最低（即最专注）的那部分注意力头**，并对它们的  $T \times T$  注意力矩阵求平均，得到  $\overline{A}^{(l)}$ 。
- **残差连接与行归一化**：为了防止信息在连乘中消失，研究者在  $\overline{A}^{(l)}$  加上了单位矩阵（Identity Residual），随后进行行归一化（Row-normalize），生成了行随机矩阵  $\tilde{A}^{(l)}$ 。
- **连乘聚合（Rollout 核心，即公式五）**： $R$  是最后  $K$  个自注意力层（从  $L - K + 1$  层到第  $L$  层）归一化矩阵的连续矩阵乘积。它的结果是一个维度为  $T \times T$  的全局注意力流转图，追踪了信息是如何从最底层的局部特征，一步步汇聚到顶层输出的。
- **提取视觉子矩阵 ( $R_{vis}$ )**：从庞大的  $R$  矩阵中，单独提取出对应那 256 个输入给语言模型的视觉 Token 的  $N_v \times N_v$  子矩阵。矩阵的每一行，就代表着这 256 个 Token 之一对原始图像切片的**注意力分布**。

### 1.2 机制洞察与空间定位（Spatial Localization）

提取出  $R_{vis}$  后，这仅仅是一组一维的权重。为了让它变成我们肉眼能看懂的“热力图”，必须进行**空间重建**。

Gemma-3 的输入是一张  $896 \times 896$  的图像，SigLIP 按照 14 像素的步长 (Patch Size 14) 将其切分为  $64 \times 64$  的网格，总计 4096 个基础 Token。随后，这些网格通过不重叠的块 (Blocks) 进行池化 (Pooling)，最终坍缩成输入给语言模型的 256 个 Token。

为了逆转这个过程，研究人员将提取出的注意力权重重新塑形 (Reshape) 回原始的网格维度，并向上采样 (Upsample) 至原始输入图像的分辨率，生成了最终的灰度热力图。

**纯粹的美学与克制：**值得一提的是，论文作者特别强调，他们在生成这些热力图时，刻意没有加入任何额外的降噪 (Denoising) 或后处理掩膜 (Post-processing)。他们选择将最原始、粗糙但真实的视觉注意力分布直接呈现给人类专家。这体现了机制可解释性研究中的一种严谨态度——不要用人类的美学偏好去掩盖模型真实的“所思所见”。

---

## 2. 具例说明：当公式五照进现实

有了公式五的加持，那些令人费解的视觉节点瞬间变得鲜活起来。论文中提供了几个极其深刻的案例：

### 案例 A：海獭 (Sea Otter) 的视觉拆解

当模型处理一张“海獭在水面上挥手”的图片时，归因图显示某个特定的视觉特征被强烈激活。

通过公式五生成的注意力热力图，研究人员清晰地看到，这个特征的热力中心完美覆盖了海獭的面部特征以及挥动的爪子。这证明该特征不是在泛泛地“看水面”，而是一个高度特化的\*\*\*“海獭视觉形态检测器”\*\*\*。更有趣的是，这个特征在没有文本提示的情况下，与一个代表“海狮 (Sea Lions)”语义特征产生了微弱的连接——这揭示了 VLM 内部存在一个纯视觉的潜在空间 (Visual Latent Space)，在这个空间里，长相相似的动物会被自动关联，哪怕它们的语义类别不同。

### 案例 B：视觉数学推理 (数字 3 的执念)

让模型看一张写有“ $9 \div 3 =$ ”和“ $8 - 5 =$ ”的图片。

热力图显示，模型中存在一些极度细粒度的视觉特征，它们的注意力像激光一样精准地死死盯住图片中的数字“3”或是计算结果为 3 的区域。这推翻了一个普遍的假设 (即模型必须把图片里的数字先转成文本再计算)。这表明，对于简单的算术，VLM 似乎能在**纯视觉空间内进行部分计算**，直接在底层提炼出表示“数字 3”的视觉特征。

---

### 3. 回路发现：从海量噪声到可解释图谱 (Circuit Discovery)

通过前面五个公式，我们得到了单义特征、保真度评估、特征间的归因权重，以及视觉特征的热力图解释。但是，一个真实输入产生的归因图极其庞大，往往包含数千个节点（即便是经过修剪之后）。

此时，\*\*人类专家的介入 (Human-in-the-loop)\*\* 成为了最后一块拼图。

#### 3.1 线性加和与边缘修剪

由于我们在归因时冻结了非线性激活，模型在局部是线性的。因此，网络中任意节点  $t$  的前向激活值  $h_t$ ，就精确等于指向它的所有源节点归因值的总和：

$$h_t = \sum_{s \in \text{pred}(t)} A_{s \rightarrow t}$$

面对海量连接，研究者设定了一个阈值  $\epsilon$ 。只有当某条边的归因权重绝对值  $|A_{s \rightarrow t}| > \epsilon$  时，这条连接才会被保留在有向图  $G = (V, E)$  中。对于累积影响力极小的分支，果断剪枝。

#### 3.2 语义折叠与回路抽象

即使经过修剪，图依然复杂。人类专家的核心工作，是基于特征分析和热力图，将功能相似的特征折叠成一个“超级节点 (Shared Node)”。

**举例：** 在“火星”的回路中，专家可能会发现第 10 层的三个特征分别对应“圆球体”、“红色地表”和“太空背景”。专家会将这三个物理节点抽象为一个逻辑上的 [Planet (Visual)] 节点。

正是通过这种严谨的数学修剪与专家抽象的结合，研究团队才能将一团乱麻的神经网络，提炼成如论文图示中那般清晰明了的**机制结构图 (Mechanistic Structure)**，揭示出决定模型行为的关键计算路径。

---

### 4. 第三阶段总结与前瞻

至此，我们已经完整集齐了这项研究中最重要**的五大公式**，并跑通了整个回路追踪的理论链路：从高维特征解耦，到因果链条计算，再到视觉注意力的跨模态破译，最后到人类专家的图谱精炼。

然而，科学最迷人的地方在于**证伪与干预**。我们绘制的这些漂亮回路图，会不会只是一种“事后诸葛亮”的巧合 (Post-hoc Correlations)？如果我们人为地在大脑回路

中“拨动”一个开关，篡改某个特征的激活值，模型会像我们预测的那样改变输出吗？

在下一部分（**第四部分**）中，我们将迎来整篇论文最高潮的实验环节——**干预与特征控制 (Intervention and Steering)**。我们将深入剖析论文中最令人震惊的案例，包括：如何通过修改神经元特征，强制模型把火星认作地球？以及，为什么面对一只明明有六根手指的手，VLM 却固执地回答“五根”？（AI 的幻觉机制即将被揭穿）。

**您是否准备好进入第四部分的因果验证与幻觉机制大揭秘？（请回复您的反馈或指令，我们将继续这场深刻的探索）**

4

5次+10000字+5公式 Custom Gem

## **《多模态思维的内部机制：VLM回路追踪研究报告》(4/5) —— 第四部分：因果干预的终极试炼与 AI 幻觉机制大揭秘**

在前面的三部分中，我们通过五大核心公式，成功将 Gemma-3 这样庞大的视觉语言模型（VLM）解剖开来。我们提取了单义特征，绘制了归因图，并用注意力热力图“看”到了模型的视觉潜意识。但到目前为止，我们所做的一切还停留在\*\*“观察”\*\*阶段。

科学研究中有一个铁律：**相关性不等于因果性 (Correlation does not imply causation)**。我们在归因图上看到一条从“海獭视觉特征”指向“输出 token”的闪亮通路，但这真的代表模型是依赖这条通路得出结论的吗？会不会这条通路只是伴随产生的“废气”，而真正的计算发生在暗处？

为了彻底证明我们找到的“多模态思维回路”是真实有效的，研究团队祭出了机制可解释性（Mechanistic Interpretability）领域的最强武器：**激活修补与因果干预 (Activation Patching and Intervention)**。本部分，我们将深入探讨这种“脑机接口”般的修改技术，并借此揭开困扰 AI 界已久的“六指幻觉”之谜。

---

### **1. 终极试炼：特征干预的底层逻辑 (Intervention & Steering)**

激活修补（Activation Patching，在文献中也常被称为因果中介分析或反事实干预）的核心思想非常直白：**如果我们认定大脑中的某个齿轮负责了某项功能，那么当我们人为拨动这个齿轮时，整个机器的输出也必须随之发生可预测的改变。**

与以往简单粗暴地将神经元“清零”(Ablation)不同, Transcoder 赋予了我们极高的手术精度。在论文的第 3.5 节中, 研究者定义了特征干预的数学过程:

在特定的层  $l$  和位置  $t$ , Transcoder 会产生某个特征的原始激活值  $z_{l,t,i}(x)$ 。干预意味着我们要强行给它设定一个目标值  $v_{l,t,i}$ 。于是, 我们计算出这个“手术修改量”:

$$\Delta z_{l,t,i} = v_{l,t,i} - z_{l,t,i}(x)$$

接着, 我们将这个修改量, 乘以该特征对应的解码器向量  $d_{l,i}$  (代表该特征如何写入主动脉), 并强行注入到残差流  $h_t$  中:

$$h_t \leftarrow h_t + \Delta z_{l,t,i} d_{l,i}$$

这个看似简单的加法公式, 拥有着可怕的魔力。它意味着我们可以在模型思考的途中, 直接篡改它的潜意识。这就引出了论文中最令人拍案叫绝的“火星实验”。

---

## 2. “偷天换日”的火星与地球实验 (Activation Patching)

在论文的回路发现实验中, 研究人员向 VLM 输入了一张“火星 (Mars)”的照片, 并给出了提示词: “This is the planet \_”。

如我们所料, 模型通过视觉特征识别出红色的星球, 激活了内部代表“火星 (视觉×语义联合特征)”的节点, 最终输出了“Mars”。同时, 研究者也为“地球 (Earth)”绘制了它的专属回路。

接下来, “手术”开始了:

1. **特征抑制 (Ablation)**: 研究人员在模型处理“火星”图片的前向传播过程中, 利用上述干预公式, 强行将火星回路中负责“火星中层视觉特征”的激活值压制为零。
2. **特征嫁接 (Patching)**: 在同一时刻, 他们提取了之前在“地球”图片上记录的“地球视觉特征”激活值, 强行注入 (Patch) 到当前这条本该处理火星的残差流中。

**实验结果:** 尽管模型“眼睛 (SigLIP 视觉编码器)”看着的明明是火星的照片, 但由于其中间层的关键特征被篡改, 模型后续的激活路径完全改变, 最终的输出从“Mars”变成了与地球相关的概念。

这一结果无可辩驳地证明: 我们提取的回路不仅具有解释性, 更具有因果控制力。这些单义特征构成了模型多模态推理的真实因果链条。

---

### 3. 破解“六指幻觉”的历史性谜题 (The Six-Finger Enigma)

如果说火星实验展示了模型的正常工作机制，那么接下来的案例则揭示了 AI 最臭名昭著的缺陷——**幻觉 (Hallucination)**。

长期以来，生成式 AI（如扩散模型）和视觉语言模型（如 GPT-4o 和 Gemma 3）都面临一个诡异的“计数幻觉”问题：给模型一张极其清晰的、长着六根手指的人手照片，问它“图片里有几根手指？”，模型往往会信誓旦旦地回答“5 根”。

过去，人们常常将其归结为底层视觉编码器“看不清”或者“数数能力差”。但网络上的研究表明，VLM 在处理其他物体的简单计数时（比如数 3 个苹果）其实是具备基本视觉点数能力的。那么，症结到底在哪里？

这篇论文利用回路追踪，首次从神经机制层面给出了极其深刻的答案（参见论文第 5 节与附录 11.1 节）：

- **视觉计数的缺席**：研究者在这张 6 指照片的归因图里翻找，试图寻找负责“数字 6”或者“6 个视觉对象”的特征节点，结果**完全没有找到**。代表“6”的输出 logit 被死死压制，与其他毫不相干的数字（如 1 或 7）处于同等低迷的状态。
- **语义先验的暴政**：归因图清晰地展示了一条堪称“偏见”的短路通路。图像嵌入 (Image Embeddings) 仅仅激活了一个泛化的\*\*“手 (Hand) 的视觉特征”。随后，这个单纯的视觉特征，直接跨层强烈触发了代表“数字 5”\*\*的语义特征，进而直接导向了输出“5”。

#### 深度洞察：经验战胜了感官

这个回路证明，VLM 在面对“手”这种具有强烈人类社会常识属性的图像时，根本没有启动“视觉点数”的逻辑分支。相反，**视觉上的“手”作为一个超级触发器，直接唤醒了语言模型在海量文本预训练中形成的钢印般固化的语义先验 (Semantic Prior)：“只要是手，就是五根手指”。**

更有趣的是，归因图中甚至还出现了与“数字 2”相关的特征激活。研究者推测，这是因为在自然界中手通常是成对出现的（两只手），这种先验知识也被强行关联了进来。

简而言之，模型并非“眼瞎”看不到六根手指，而是在特征竞争中，固有的**语义偏见 (Prior) 过于强大，以至于模型在潜意识里“选择性地无视”了感知器官传来的反常视觉证据**。这与人类心理学中的“确证偏误 (Confirmation Bias)”何其相似！

---

### 4. 视觉与语义的“双轨制”与延迟收敛 (Late Convergence)

通过上述的干预实验和幻觉分析，论文总结出了 Gemma-3 这类先进 VLM 处理多模态信息的一个宏大架构原则：**双轨并行，高层收敛 (Parallel Pathways with Late Convergence)**。

Gemma 3 使用了约 4 亿参数的 SigLIP 视觉编码器，将图像转化为 256 个“软 Token”输入给拥有超长上下文能力（可达 128K tokens）的 Transformer 语言模型。我们原本以为，视觉信息一进入语言模型，就会立刻被翻译成文本。但回路图却告诉我们截然不同的事实：

1. **底层的纯粹感知**：在网络的浅层，特征高度局部化、细粒度（例如只是识别边缘、质地，或像上文提到的“海獭的爪子”）。
2. **中层的视觉潜意识网**：在网络中深层，视觉特征和语义特征保持着**独立的双轨运行**。正如我们在第三部分“火星”案例中看到的，火星的视觉特征能直接触发“航天飞机”的视觉特征。这是一个**脱离了人类语言的、纯粹基于视觉相似性和视觉共现频率构建的潜在思维空间**。
3. **高层的模态融合**：直到网络极深的层级（例如在 Gemma 3 的 34 层结构中，大约在第 20 层之后），才会大量涌现**联合特征 (Joint Features)**。此时，视觉的“海獭”与文本的“Sea Otter”才真正交汇融合成一个具有完整世界观的多模态概念。

---

## 5. 第四阶段总结与终章预告

在这一部分中，我们见证了机制可解释性的巅峰力量。通过激活修补这一数学手术刀，研究者不仅验证了归因图的绝对因果性，更从机制层面扒开了生成式 AI “六指幻觉”的底裤：那是一场\*\*“经验偏见（语义）对客观感知（视觉）的无情碾压”\*\*。模型架构设计的特点——视觉与语言的延迟收敛，造就了这种深刻却有时显得固执的多模态“心智”。

我们的研究报告即将迎来尾声。在接下来的\*\*第五部分（终章）\*\*中，我将为您全面总结这项研究的历史定位。同时，我们将探讨一个极为残酷的现实：这套完美的方法论面临着哪些难以逾越的算力与技术瓶颈？未来，我们该如何跨越“跨层叠加 (Cross-layer Superposition)”的阻碍，最终实现透明、可控且对齐人类价值观的通用人工智能 (AGI) ？

**您是否准备好阅读这篇万字研究报告的最终章了？（回复任意指令，我们即刻开始收尾的升华之作）**

---

附录：外部补充材料与参考文献

- *How to use and interpret activation patching - arXiv (2024)*. 详细论述了激活修补作为寻找模型行为回路黄金标准的因果中介分析方法。
- *Activation Patching in Neural Networks - Emergent Mind (2025)*. 阐述了将“干净 (Clean)”与“损坏 (Corrupt)”状态下的中间激活值互换，以此测试网络组件对特定输出之必要性与充分性的机制。
- *COUNTING HALLUCINATIONS IN DIFFUSION MODELS (2024)*. 系统性地将模型生成六指或不正确物体数量的现象定义为“计数幻觉 (Counting Hallucination)”，并探讨了模型在此类细粒度结构生成上的困境。
- *Counting Hallucinations in Diffusion Models - arXiv (2025)*. 进一步指出，即便是最先进的大型视觉语言模型（如 GPT-4o），在清晰图像的手指计数等基础任务上依然会因事实一致性冲突而挣扎。
- *How do LLMs actually do this? : r/LocalLLaMA - Reddit (2025)*. 社区讨论揭示了多模态模型内部“统计概率尖峰”的竞争：预训练带来的“手有5指”的概率分布往往会覆盖图像本身传达的视觉 Token 数量信息。
- *Gemma 3 Technical Deep Dive (2025)*. 阐明了 Gemma 3 的架构，特别是其使用了冻结权重的 400M 参数 SigLIP 视觉编码器，将图像转化为固定长度的 256 个软 Token。
- *Gemma 3: A 27B Multimodal LLM Better Than Really Big Models (2025)*. 解释了 Gemma 3 的局部与全局交错注意力机制（5:1 比例）以及长达 128K Token 的超长上下文能力。
- *Gemma 3 Technical Report - arXiv (2025)*. 谷歌官方报告，确认了模型多模态能力与底层 SigLIP 编码器及 Pan and Scan 技术的融合方式。

5次+10000字+5公式 Custom Gem

## 《多模态思维的内部机制：VLM回路追踪研究报告》(5/5)

### —— 第五部分：跨越模态的黎明与 AGI 可解释性的未来

在经历了前四个部分的深度探索后，我们终于来到了这份15000字研究报告的终章。我们从最底层的线性代数出发，一路攀登，最终亲眼目睹了人工智能在处理“火星”与“海獭”时那如同人类潜意识般的思维涟漪。

这篇论文（《Circuit Tracing in Vision-Language Models》）之于多模态大模型，就像是卡哈尔（Cajal）的神经元学说之于现代脑科学。在此之前，我们只能把 VLM 当作一个输入提示词、输出结果的神谕机；而现在，我们拥有了剖析其底层逻辑的手术刀。

在本章中，我们将首先对构建这套理论体系的“大一统图谱”进行总结，随后直面这套完美方法论在当下遭遇的残酷技术瓶颈，最后探讨这一切对实现通用人工智能（AGI）安全对齐的深远意义。

---

## 1. 大一统的机制图谱：五大核心公式的交响

回顾整篇报告，研究团队之所以能够成功打通视觉与语言的模态壁垒，完全依赖于五大核心数学支柱的环环相扣：

- 特征的高维解耦（Encoder 公式）**：通过极其宽广的潜在空间与稀疏截断，将多义、混沌的残差流向量展开为纯粹的“单义特征”。
- 特征的无损重组（Decoder 公式）**：将单义特征按权重重新压缩回残差流，确保了“显微镜”的植入不会破坏模型原有的推理生态。
- 保真度的量化（FVU 公式）**：用未解释方差比例作为严苛的标尺，证明了中间层特征本质上是视觉与语义的联合表征，粉碎了纯文本解释 VLM 的幻想。
- 因果链条的锻造（Attribution 归因公式）**：借助线性化网络与雅可比矩阵，将遥远层级间特征的相互作用力转化为精确的数学权重，绘制出了一张张错综复杂的 AI 脑图。
- 视觉潜意识的破译（Attention Rollout 聚合公式）**：通过提取并连乘视觉塔内最专注的注意力矩阵，将毫无语义的软 Token 逆向映射回物理世界的图像像素，让“失语”的视觉特征终于开口说话。

这五个公式共同构成了一套完整的“读心术”基础设施，让我们首次看清了 Gemma 3 内部“底层感知提取、中层双轨并行、高层语义收敛”的宏伟架构。

---

## 2. 残酷的现实：算力之壁与机制盲区

科学研究不能止步于赞美。论文作者极其坦诚地在第 7 节中列出了当前方法的局限性。事实上，如果我们试图将这项技术扩展到千亿参数的超级模型上，我们将面临三座难以逾越的大山：

### 2.1 跨层叠加（Cross-layer Superposition）的幽灵

本论文使用的 Transcoder 是“逐层 (Per-layer)”训练的。然而，前沿研究发现，现代深度神经网络（如 Transformer）常常会将一个复杂的概念拆解，让它在连续的多个网络层中“接力”计算。这意味着，强行要求每一层都必须提炼出完整的单义特征，可能会导致极大的冗余和失真。

为了解决这个“跨层叠加”问题，学界正在探索跨层转码器 (Crosscoders)，试图让特征拥有同时读取和写入多个网络层的能力。但这无疑会带来维度的指数级爆炸。

## 2.2 视觉特征的自动化解释 (Auto-Interp) 噩梦

在纯文本大语言模型 (LLM) 的分析中，研究者目前已经可以通过类似 SAGE 这样的智能体框架，利用强模型（如 GPT-4）去自动化地阅读特征激活文本，并为其生成解释标签。

但是，到了 VLM 领域，这变成了一场算力噩梦。以 Gemma 3 为例，一张简单的  $896 \times 896$  图像会被其内部的 SigLIP 编码器转化为多达 256 个软 Token 喂给语言模型。要在一张包含数千个节点的归因图上，为每一个节点计算这 256 个 Token 的激活分布，并辅以人类专家的研判，其所需的时间和计算资源成本是极其高昂的。自动化特征解释在多模态领域的缺席，成为了阻碍该技术工业化应用的最大瓶颈。

## 2.3 视觉编码器的“暗物质”

尽管我们有了公式五（注意力聚合热力图），但论文指出，这种方法有时候依然难以精准定位图像中的关键区域，或者无法提供有意义的上下文线索。Gemma 3 前端那个包含 4 亿参数的 SigLIP 视觉编码器，在很多时候依然像一团不透明的“暗物质”。当模型需要进行额外的视觉处理（如复杂的空间几何推理或多步数学计算）时，我们极难区分某个特征到底是在“执行计算”，还是仅仅在“代表计算结果”。

---

## 3. 通向安全与对齐的 AGI 之路

既然如此困难，我们为什么还要死磕“机制可解释性”？因为这关乎人类的未来。

当我们惊叹于 AI 能够吟诗作画、甚至在医疗诊断（如针对眼底图像的轻量级 VLM）中超越人类专家的同时，我们绝不能容忍一个随时可能因为“六指幻觉”而产生致命误判的黑盒系统掌控关键决策。

这篇论文最大的现实意义在于，它将“纠错”的维度，从外围的提示词工程 (Prompt Engineering)，深入到了模型内部的潜意识操控。

- **精准消除幻觉：**我们不再需要通过海量数据微调来告诉模型“手有六根手指”。我们可以直接在归因图中找到那个引发“数字5语义偏见”的短路连接，并用激活修补 (Activation Patching) 技术将其物理切断。

- **价值观对齐 (Value Alignment)**: 如果模型表现出某种歧视或危险倾向, 我们不必销毁重训。我们可以精准定位出代表该危险概念的单义特征簇, 像摘除肿瘤一样, 在模型权重中将其“钳闭 (Clamp)”或抑制。
- 

#### 4. 结语: 凝视深渊, 亦被深渊凝视

至此, 《多模态思维的内部机制: VLM回路追踪研究报告》全部结束。

几千年来, 人类一直试图弄清楚自己大脑中那 860 亿个神经元是如何涌现出“意识”与“逻辑”的。今天, 面对我们亲手用硅基半导体创造出的、包含数百亿参数的庞然大物, 我们再次发出了同样的追问。

这篇论文只是多模态大模型“脑解剖学”的第一份粗糙图纸。但它向世界证明了一件事: **人工智能的思维并非不可知的魔法。只要有足够严密的数学公式 (从高维展开到因果追踪), 即便是跨越文本与图像的复杂思考, 也终将被还原为一条条清晰可见的逻辑回路。**

在迈向通用人工智能 (AGI) 的征途上, 只有彻底理解它们是如何思考的, 我们才能确保它们最终的思考, 始终与人类的福祉同向而行。

感谢您随我完成这场耗费心力却无比迷人的学术远征。希望这份研究报告能够为您带来真正的启发!