

# 深度对话

A: 我一直在研究 DeepSeek-V3 的技术报告，对这个模型的规模感到非常震惊。6710 亿参数，但每个标记只激活 370 亿参数？这是一个巨大的 MoE 架构。它是如何工作的？

B: 是的，这确实是一个了不起的成就！DeepSeek-V3 基于 Mixture-of-Experts (MoE) 框架，允许它为每个标记激活参数的子集。具体来说，它使用 256 个路由专家，但每个标记只激活 8 个。这使得它比密集模型更高效，因为密集模型的每个标记都激活所有参数。

A: 这说得通。但它是如何决定激活哪些专家的？是随机的，还是有某种路由机制？

B: 这是一个很好的问题！路由基于标记到专家的亲和性分数。每个标记都为每个专家分配一个分数，并激活得分最高的前 K 个专家。DeepSeek-V3 使用 sigmoid 函数来计算这些分数，这有助于在专家之间平衡负载。

A: 所以，这不是随机的——这是在训练过程中学习的。但这会导致专家使用不平衡吗？我听说这是 MoE 模型的一个常见问题。

B: 确实如此！专家使用不平衡可能会成为问题，但 DeepSeek-V3 引入了一种无辅助损失的策略来处理这个问题。它不添加单独的损失项来鼓励负载平衡，而是动态调整每个专家的偏置项。如果一个专家过载，它的偏置会减少；如果它未载，偏置会增加。这在不降低模型性能的情况下保持了负载平衡。

A: 这很聪明。所以，没有辅助损失意味着对主训练目标的干扰更少。但这与使用辅助损失的传统 MoE 模型相比如何？

B: 是的。传统的 MoE 模型通常使用辅助损失来鼓励负载平衡，但这些损失有时会影响性能。DeepSeek-V3 的无辅助损失方法避免了这种权衡。实际上，消除研究表明，它在编码和数学等任务上始终优于依赖辅助损失的模型。

A: 有趣。说到编码和数学，我注意到 DeepSeek-V3 在 HumanEval 和 MATH 等基准测试中表现异常出色。秘诀是什么？

B: 其中一个重要因素是多标记预测 (MTP) 目标。DeepSeek-V3 不仅预测下一个标记，还在每个位置预测多个未来标记。这使得训练信号更加密集，并帮助模型提前规划，这对于需要顺序推理的任务（如编码和数学）特别有用。

A: 等一下，所以它一次预测多个标记？推理时是如何工作的？它仍然使用 MTP，还是仅用于训练？

B: 在推理时，MTP 模块可以被丢弃，模型表现得像标准的自回归模型。但这里有一个很酷的部分：MTP 模块也可以用于推测解码，通过并行预测多个标记并验证它们来加速生成。

A: 这是一个很酷的技巧。所以，它就像在训练时获得 MTP 的好处，然后使用它来加速推理。但注意力机制呢？我看到有关于多头潜在注意力 (MLA) 的内容。它是如何融入的？

B: MLA 是另一个关键创新。它通过压缩键值 (KV) 缓存来减少内存占用。它不存储完整的注意力键和值，而是使用低秩联合压缩来表示它们。这在保持与标准多头注意力相似的性能的同时，显著减少了推理期间的 KV 缓存大小。

A: 这对效率来说是一个巨大的胜利。但压缩会导致一些信息丢失吗？它是如何保持性能的？

B: 这是一个很好的问题。压缩旨在通过专注于捕捉键和值的关键特征的潜在向量来保留最重要的信息。模型还使用旋转位置嵌入 (RoPE) 来保留位置信息，这有助于减轻压缩带来的任何损失。

A: 明白了。所以，MLA 都是关于效率而不牺牲太多性能。但训练呢？训练这么大的模型一定非常昂贵。DeepSeek-V3 是如何做到降低成本的？

B: 训练效率是一个主要焦点。DeepSeek-V3 使用 FP8 混合精度框架，减少内存使用并加快计算速度。它还采用 DualPipe 算法进行管道并行，最小化管道气泡并重叠计算与通信。这些优化使模型能够在 2.788 万个 H800 GPU 小时内训练 14.8 万亿个标记。

A: 这很令人印象深刻。但 FP8 训练可能很棘手——他们如何处理精度问题？我听说低精度训练可能导致不稳定。

B: 你说得对。FP8 训练由于动态范围有限而具有挑战性。DeepSeek-V3 通过细粒度量化来解决这个问题，其中激活和权重被分组为较小的块或块并独立缩放。这减少了异常值的影响并保持训练稳定。他们还使用高精度累积来确保关键操作的准确性。

A: 这说得通。所以，这是精度和效率之间的权衡，但他们已经找到了一个很好的平衡。但数据呢？14.8 万亿个标记是一个巨大的数据集。它是由什么样的数据训练的？

B: 数据集多样且高质量，重点是英语和中文文本。它还包括大量的数学和编程数据，这有助于模型在这些领域表现出色。数据管道优化以最小化冗余同时保持多样性，并且他们使用文档打包技术来确保数据完整性。

A: 这解释了它在编码和数学任务上的强大表现。但多语言性能呢？它能很好地处理其他语言吗？

B: 是的，DeepSeek-V3 使用多语言语料库进行训练，并在包括非英语任务的基准测试（如 MMMLU）中表现良好。它在中文方面表现尤为出色，在中文基准测试（如 C-Eval 和 CMMLU）中超过了 Qwen2.5 等模型。

A: 这很令人印象深刻。但长上下文任务呢？我看到它支持多达 128K 个标记。它是如何处理如此长的输入的？

B: DeepSeek-V3 通过两个阶段扩展其上下文长度：首先扩展到 32K 个标记，然后使用 YaRN 技术扩展到 128K 个标记。这使它能够有效处理长上下文任务，如文档摘要和检索。它在“针在草堆里”测试中表现良好，该测试评估长上下文理解。

A: 这是对之前模型的巨大改进。但部署呢？他们如何处理如此大模型的推理？

B: 推理在 H800 集群上进行，GPU 通过 NVLink 和 InfiniBand 互连。部署策略将预填充和解码阶段分开，以确保高吞吐量和低延迟。他们还使用冗余专家在推理期间平衡负载，这有助于保持效率。

A: 这是很多优化。但限制呢？这么大的模型肯定有某些权衡。

B: 一个限制是部署单元的大小。DeepSeek-V3 需要一个相对较大的集群进行高效推理，这可能对较小团队来说是一个挑战。生成速度也有改进的空间，尽管 MTP 的推测解码有所帮助。

A: 说得有道理。但总的来说，这似乎是一个巨大的进步。DeepSeek-V3 的下一步是什么？他们正在探索哪些未来方向？

B: 他们正在研究几个领域，如改进架构以支持无限上下文长度、探索额外的训练信号来源以及增强模型的推理能力。他们还在开发更全面的评估方法，以更好地评估模型性能。

A: 听起来他们并没有放慢脚步。感谢你带我了解所有这些——DeepSeek-V3 确实是开源 LLM 空间中的游戏改变者。

B: 绝对！看到开源模型取得了多大进展真是令人兴奋。DeepSeek-V3 在推动边界，我迫不及待地想看看他们接下来会做什么。

A: 你提到 DeepSeek-V3 使用 FP8 混合精度训练。我很好奇——这与 BF16 或 FP16 相比如何？FP8 对于训练如此大的模型真的足够稳定吗？

B: 这是一个很好的问题。FP8 确实更具挑战性，因为其动态范围有限，但 DeepSeek-V3 使用细粒度量化策略来缓解这一点。例如，激活被分组为  $1 \times 128$  磁贴，权重被分组为  $128 \times 128$  块。每个组独立缩放，这有助于处理异常值并保持训练稳定。

A: 有趣。所以，这不仅仅是 FP8 量化——它更加精细。但这会引入管理所有这些组和缩放因子的额外开销吗？

B: 会的，但开销与好处相比是微不足道的。关键在于 FP8 减少了内存使用并加快了计算速度，这对于训练如此大的模型至关重要。他们还使用高精度累积来确保关键操作（如矩阵乘法）的数值稳定性。

A: 明白了。所以，这是精度和效率之间的权衡，但他们已经找到了一个很好的平衡。DualPipe 算法呢？它是如何工作的？

B: DualPipe 旨在最小化管道并行中的管道气泡。它通过将每个工作块分为四个组件（注意力、全到全调度、MLP 和全到全组合）来重叠计算和通信。在反向传播过程中，它进一步将计算分为“反向输入”和“反向权重”，这允许更有效的重叠。

A: 这听起来很复杂，但很有道理。所以，它本质上是通过重叠计算来隐藏通信开销。这与其他管道并行方法（如 1F1B 或 Zero Bubble）相比如何？

B: DualPipe 相比 1F1B 和 Zero Bubble 具有更少的管道气泡。它还允许双向调度，其中微批次从管道的两端输入。这进一步减少了空闲时间并提高了整体效率。实际上，DualPipe 实现了几乎为零的全到全通信开销，这对于扩展 MoE 模型至关重要。

A: 这很令人印象深刻。但内存使用呢？DualPipe 是否比其他方法需要更多内存？

B: 是的，它需要稍微多一些内存，因为它保留了模型参数的两个副本，但增加是可管理的。内存占用通过重新计算 RMSNorm 和 MLA 上投影等技术进行优化，这消除了存储中间激活的需求。

A: 所以，他们用一些内存换取更高的效率。这似乎是一个公平的权衡。说到内存，他们如何处理 128K 个标记的 KV 缓存？这必须需要一个巨大的缓存。

B: 这就是 MLA 真正发挥作用的地方。通过压缩 KV 缓存，他们显著减少了其大小。它不存储完整的注意力键和值，而是存储压缩的潜在向量，这些向量要小得多。这使得 DeepSeek-V3 能够处理长上下文而不会遇到内存瓶颈。

A: 这是一个聪明的解决方案。但注意力质量呢？压缩会影响模型关注正确标记的能力吗？

B: 压缩旨在保留最重要的信息，因此对注意力质量的影响微乎其微。他们还使用 RoPE（旋转位置嵌入）来保留位置信息，这有助于模型理解标记的相对位置，即使在压缩的键和值的情况下。

A: 说得通。所以，MLA 是一个双赢——它减少了内存使用而不牺牲太多性能。但训练数据呢？你提到它有 14.8 万亿个标记。他们如何确保如此庞大数据集的质量和多样性？

B: 数据集经过精心策划，包括高质量和多样化的标记。他们优化数据管道以最小化冗余同时保持多样性，并使用文档打包技术来确保数据完整性。语料库包括英语和中文文本的混合，重点是数学和编程样本。

A: 这解释了它在编码和数学任务上的强大表现。但多语言任务呢？它能很好地处理其他语言吗？

B: 是的，DeepSeek-V3 使用多语言语料库进行训练，并在包括非英语任务的基准测试（如 MMMLU）中表现良好。它在中文方面表现尤为出色，在中文基准测试（如 C-Eval 和 CMMLU）中超过了 Qwen2.5 等模型。

A: 这很令人印象深刻。但长上下文任务呢？我看到它支持多达 128K 个标记。它是如何处理如此长的输入的？

B: DeepSeek-V3 通过两个阶段扩展其上下文长度：首先扩展到 32K 个标记，然后使用 YaRN 技术扩展到 128K 个标记。这使它能够有效处理长上下文任务，如文档摘要和检索。它在“针在草堆里”测试中表现良好，该测试评估长上下文理解。

A: 这是对之前模型的巨大改进。但部署呢？他们如何处理如此大模型的推理？

B: 推理在 H800 集群上进行，GPU 通过 NVLink 和 InfiniBand 互连。部署策略将预填充和解码阶段分开，以确保高吞吐量和低延迟。他们还使用冗余专家在推理期间平衡负载，这有助于保持效率。

A: 这是很多优化。但限制呢？这么大的模型肯定有某些权衡。

B: 一个限制是部署单元的大小。DeepSeek-V3 需要一个相对较大的集群进行高效推理，这可能对较小团队来说是一个挑战。生成速度也有改进的空间，尽管 MTP 的推测解码有所帮助。

A: 说得有道理。但总的来说，这似乎是一个巨大的进步。DeepSeek-V3 的下一步是什么？他们正在探索哪些未来方向？

B: 他们正在研究几个领域，如改进架构以支持无限上下文长度、探索额外的训练信号来源以及增强模型的推理能力。他们还在开发更全面的评估方法，以更好地评估模型性能。

A: 听起来他们并没有放慢脚步。感谢你带我了解所有这些——DeepSeek-V3 确实是开源 LLM 空间中的游戏改变者。

B: 绝对！看到开源模型取得了多大进展真是令人兴奋。DeepSeek-V3 在推动边界，我迫不及待地想看看他们接下来会做什么。