



Foto: Fernando Moreira / Freepik IA

A Era da Precisão: Por que o futuro da IA na engenharia de mídia é vertical, especializado e proprietário

A engenharia de mídia entrou numa etapa onde os modelos genéricos de IA já não atendem às exigências técnicas do broadcast. O futuro é vertical, com SLMs especializados operando on-premise, garantindo latência mínima, segurança e aderência absoluta às normas.

Por Vinicius Gholmie

Nos últimos dois anos, a indústria de mídia e entretenimento foi varrida por uma onda de deslumbramento tecnológico. A narrativa predominante era clara: “maior é melhor”. Fomos levados a acreditar que a solução para tudo — da automação do master control à criação de roteiros — residia em uma única superinteligência na nuvem, alimentada por LLMs (*Large Language Models*) com trilhões de parâmetros.

No entanto, à medida que a poeira do hype assenta e entramos na fase crítica de implementação nos fluxos de trabalho de engenharia, uma realidade técnica se impõe: modelos generalistas não são adequados para a precisão cirúrgica que o broadcast exige.

Para o engenheiro de TV, o futuro não está em modelos que “sabem um pouco de tudo”, mas em arquiteturas verticais que dominam profundamente o seu negócio. Estamos migrando da era da IA Genérica para a IA Vertical, Especializada e Proprietária.



Foto: Canva IA

O problema da engenharia: A “Maldição da Generalidade”

Para entender por que o modelo gigante falha na operação diária de uma TV, precisamos olhar para a arquitetura de dados. LLMs generalistas (como o GPT-4 ou Gemini) são treinados na “mídia da internet”. Eles sofrem do que na engenharia de dados chamamos de “Maldição da Dimensionalidade”.

Ao tentar abranger tudo — de física quântica a receitas de bolo — esses modelos introduzem um ruído estatístico imenso. Quando aplicados a tarefas que exigem raciocínio lógico estrito ou adesão a normas técnicas de vídeo (como classificar um padrão de compressão ou metadados de arquivo), eles falham.

Isso ocorre devido ao fenômeno da “Escala Inversa” (*Inverse Scaling*): em tarefas de alta especificidade, modelos maiores tendem a performar pior do que modelos menores. O modelo gigante é “viciado” em seus conhecimentos prévios da web aberta e tende a ignorar instruções específicas do prompt em favor da resposta estatisticamente mais provável, gerando alucinações.

O Risco Operacional é fatal para uma emissora. Um modelo treinado na internet aberta pode não distinguir a nuance semântica entre um conteúdo apropriado para uma certa audiência de outro conteúdo impróprio.

A solução técnica: SLMs e Edge Computing

A resposta para esse gargalo de eficiência são os SLMs (*Small Language Models*). Diferente dos gigantes com trilhões de parâmetros, os SLMs (como a família Phi, Gemma ou Llama-3-8b) operam na casa de 1 a 10 bilhões de parâmetros. Para a engenharia de canais, as vantagens dos SLMs são “brutais” e imediatas:

Latência e inferência local: Em operações ao vivo (live production), não podemos esperar 3 segundos por uma resposta de API na nuvem. SLMs são leves o suficiente para rodar on-premise

ou em arquiteturas de borda (Edge Computing), diretamente nas ilhas de edição ou unidades móveis, garantindo respostas em milissegundos.

Custo de computação: Rodar um modelo especializado custa uma fração do preço de inferência de um modelo generalista, viabilizando a escala econômica

Privacidade e segurança: Ao rodar localmente, os dados sensíveis da emissora não precisam trafegar para a nuvem pública, mantendo a soberania do conteúdo.

Aplicações reais: Do arquivo ao vivo

A transição para a IA Vertical já está transformando fluxos de trabalho em grandes players globais, resolvendo problemas que modelos genéricos não conseguiram tocar.

1. Gestão de Ativos e o Fim do “Dark Data”:

O maior passivo das emissoras hoje é o “Dark Data” — petabytes de conteúdo histórico invisível aos sistemas de busca. O problema do modelo genérico: Ao analisar um vídeo de arquivo, um LLM descreve a cena genericamente como “homens jogando bola”. A solução vertical: Um canal de TV, por exemplo, utilizou modelos especializados para gerar metadados quadro a quadro. Um SLM treinado nas regras do esporte e no elenco de jogadores identifica: “Atleta X comete falta tática no minuto 35”.

Resultado: Economia de mais de 6.000 horas de trabalho humano anuais e transformação de arquivo morto em ativo líquido para canais FAST ou VOD

2. Esportes Ao Vivo: Otimização de Highlights:

A latência é inimiga do engajamento. Players como

a canais de esporte multinacionais implementaram arquiteturas de inteligência vertical que utilizam sinais multimodais — visão computacional (rastreamento de bola e jogadores) combinada com a análise de áudio da vibração da torcida. **Resultado:** O tempo de criação de highlights foi reduzido de 45 minutos (processo manual) para apenas 5 minutos. Isso define quem lidera a conversa na “segunda tela”.

3. Jornalismo e “Data Provenance”: No jornalismo, a credibilidade é a moeda, e a “caixa preta” dos LLMs é um risco inaceitável. A tendência, adotada por alguns grupos de notícias, é a criação de “Jardins Murados”. Um canal de notícias britânico desenvolveu o “*Style Assist*”, um modelo treinado exclusivamente em dezenas de milhares de artigos próprios. Ele conhece o guia de estilo e os padrões éticos da casa. Isso garante a Procedência de Dados (*Data Provenance*): o modelo não alucina; ele recupera fatos validados pelo próprio acervo da emissora, funcionando como uma ferramenta de inteligência editorial auditável.

Conclusão: Não terceirize o seu cérebro

Para os executivos e engenheiros da nossa indústria, a lição para os próximos anos é estratégica: seus dados proprietários são seu maior fosso defensivo (*moat*). O valor real não está em alugar uma inteligência genérica via API, mas em capturar, organizar e processar seus próprios

dados através de modelos especializados. O futuro pertence à IA que não tenta saber “tudo sobre o mundo”, mas que sabe “tudo sobre o seu negócio”.

A era da IA genérica e curiosa acabou. Entramos na era da IA especialista, produtiva e, acima de tudo, proprietária.



Vinicius Gholmi

é empreendedor entusiasta, amante de tecnologia, fanático por esportes e entretenimento. Após mais de 15 anos no mercado financeiro, gestão de fundos de ações/derivativos, governança corporativa/empresas familiares e alguns investimentos em startups, dedicou-se a desenvolver novas tecnologias que podem potencializar a poderosa alavanca da transformação social e aprimorar as relações com a sociedade. Bacharel em Economia pelo Insper, São Paulo, pós-graduado em Derivativos pela CBOT/CME Group, Chicago, e em Mercados Comportamentais e Adaptativos pela MIT Sloan School, Boston. Investidor de startups e, desde 2013, com experiência em governança, gestão empresarial, investimentos e inovação.

Contato: vinicius@mantis-ai.com