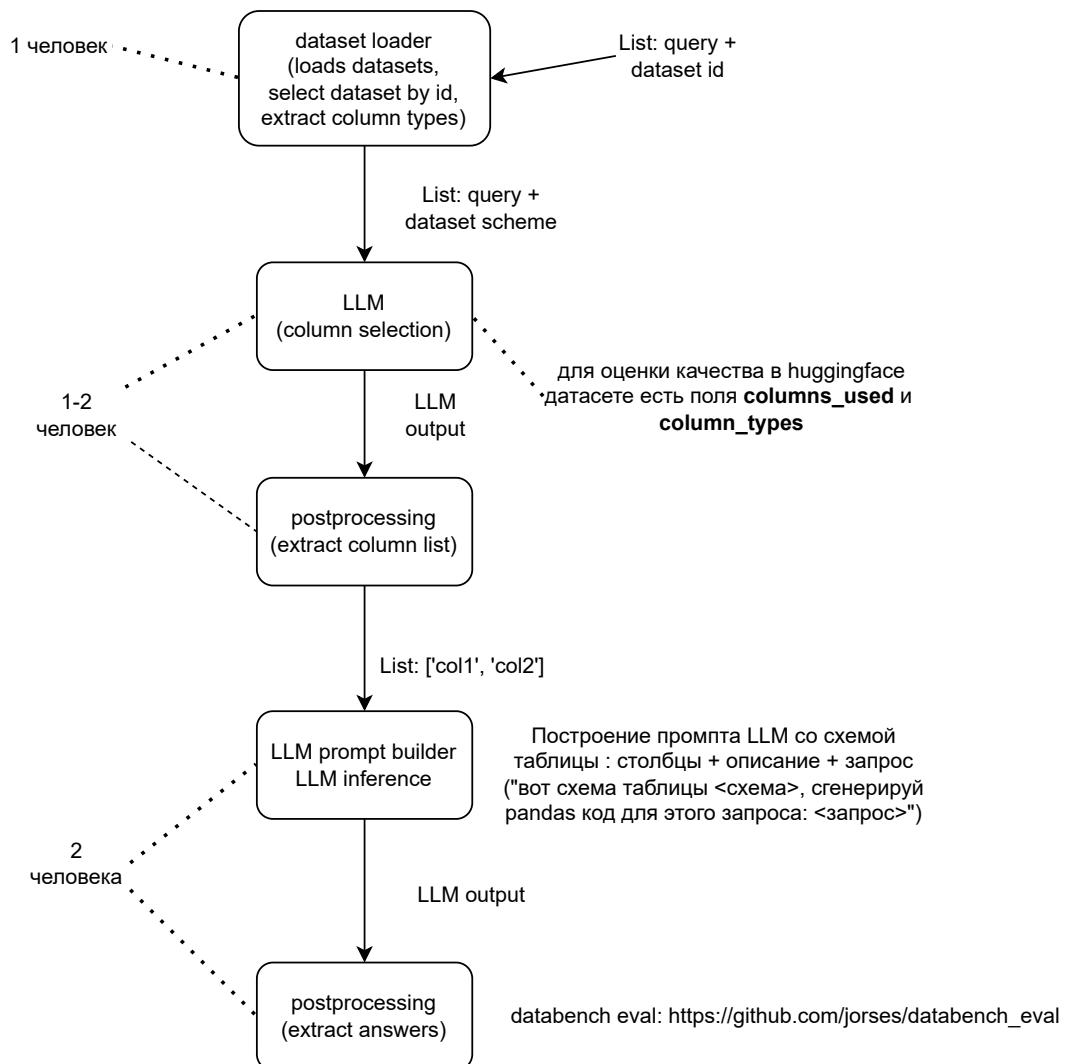
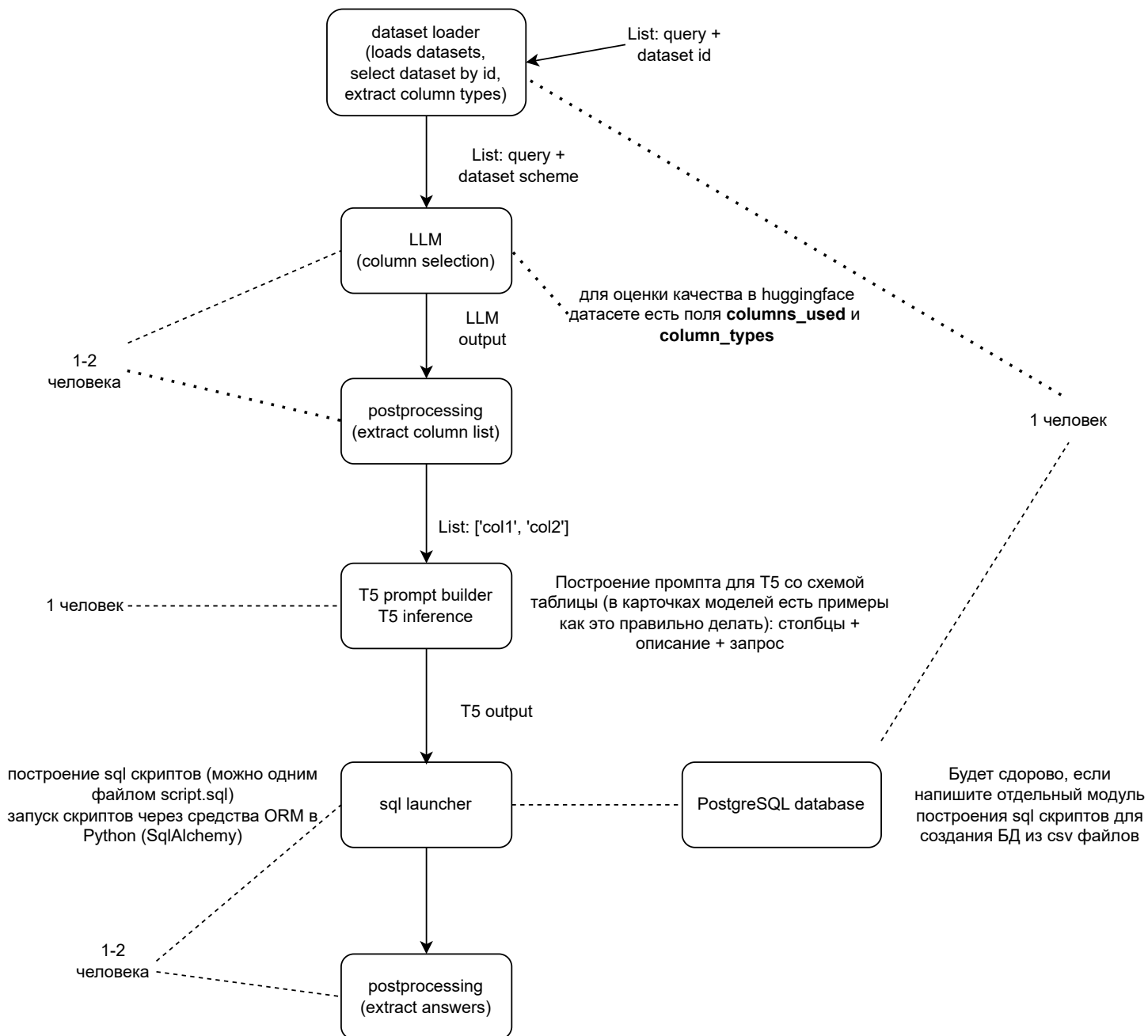


LLM



T5



LLM (column selection)	API: Llama 70B Qwen3-258B Local: DeepSeek-R1-distill-Qwen-32B, DeepSeek-R1-distill-Qwen-32B-Quant-w4a16
LLM inference (code)	Local: DeepSeek-R1-distill-Llama-70B-Quant-w4a16 Qwen2.5-Coder-32B-Instruct

ссылки на модельки (Все нормально разворачиваются через vllm):

R1 distill Llama 70B int4: <https://huggingface.co/RedHatAI/DeepSeek-R1-Distill-Llama-70B-quantized.w4a16>

R1 distill Qwen 32B: <https://huggingface.co/deepseek-ai/DeepSeek-R1-Distill-Qwen-32B>

R1 distill Qwen 32B int4: <https://huggingface.co/RedHatAI/DeepSeek-R1-Distill-Qwen-32B-quantized.w4a16>

Qwen 2.5 32B Coder: <https://huggingface.co/Qwen/Qwen2.5-Coder-32B-Instruct>

API (Qwen3, Llama 70B): <https://yandex.cloud/ru/docs/foundation-models/concepts/yandexgpt/models>