

ViT Training 2.0T -> 0.1T data CoCa-loss with tiny language decoder -> <u>align to LLM</u>

Joint Pre-training 1.4T data Up to 40% Multimodal Data

Progressive Multimodal Ratio

1 resumes LR scheduler

Joint Cooldown 0.6T data High-quality Text & Multimodal Data

Re-warmup to higher LR

RoPE base: 50.000 -> 800.000

nesumes LR scheduler

Long Text & Long Video & Long Doc

Joint Long-context

0.3T data