

Resumo: piloto migração híbrida A100 self-host

Proposta: piloto de 4 semanas migrando os 2 workflows mais caros do hub (Face Swap + FLUX dev) de fal.ai pra GPU A100 self-hosted com ComfyUI backend. Investimento R\$ 5k. Decisão de expandir baseada em gates objetivos medidos no fim do piloto.

Estado atual

Hub knight-s-forge em produção com 15 workflows (Workflow Studio JSON-driven via Supabase Edge + Postgres). Distribuição de providers:

- **fal.ai** (pay-per-request): 7 workflows críticos — Face Swap, FLUX, Lipsync v2/v3, Voice Clone, Motion Control Kling. Representa 70-80% do custo variável.
- **useapi.net** (bridge \$15/mês fixo): VEO 3 Labs, Hailuo, Midjourney, Runway, Mureka, Kling oficial, Dreamina — via contas pessoais com custo marginal \$0.
- **RunPod Serverless:** Lipsync Básico (container OSS próprio, 89% margem).
- **Supabase Pro + Vercel + Hetzner bridge:** infra base.

Hardening recém aplicado (PR #54, #55): ErrorBoundary, Sentry stub, rate-limit RPC atômico, pg_cron account health + ledger reconciliation, admin pages, Motion Control migrado pra Kling 2.1/2.6 Pro.

Custo atual: R\$ 800 - 4.500/mês (varia com volume fal.ai).

Proposta do piloto

Migrar APENAS 2 workflows pra A100 self-host:

Workflow	Provider hoje	Custo/run hoje	Provider proposto	Custo/run novo
Face Swap	fal-ai/pixverse/swap	R\$ 1.00	A100: InsightFace + InSwapper + GFPGAN + CodeFormer + Real-ESRGAN	R\$ 0.02
Imagem FLUX dev	fal-ai/flux/dev	R\$ 0.12	A100: FLUX dev 28 steps via ComfyUI	R\$ 0.01

NÃO entra no escopo: Lipsync, Voice Clone, Motion Control, Vídeo Hollywood, Música, reescrita do workflow-engine, troca de Supabase. Fica intocado.

Por que esses 2: maior volume x maior custo unitário x OSS atinge nota 10 documentado x pipeline relativamente simples x independente (não bloqueia outros workflows).

Arquitetura

Adiciona um novo case `provider: "comfyui"` no `workflow-engine.ts` existente. Não reescreve nada.

```
USER → hub frontend (Vercel)
  → POST /functions/v1/workflow-run (Supabase Edge)
    → consume_credits() atômico
    → dispatch por provider:
      → "fal" (fal.ai HTTP - mantém)
      → "useapi" (useapi.net HTTP - mantém)
      → "runpod" (RunPod serverless - mantém)
      → "comfyui" (NOVO: A100 cloud + ComfyUI + Cloudflare Tunnel)
```

Workflows novos vão como slugs paralelos (`face-swap-self` , `image-flux-dev-self`). Feature flag controla A/B test gradual.

GPU recomendada: **A100 80GB RunPod Community Cloud** (\$0.79/h ≈ R\$ 2.842/mês 24/7). Familiaridade com provider que já roda Lipsync.

Custos

Upfront (4 semanas piloto)

- Dev time: ~160h (4 semanas × 40h)
- A100 80GB 30 dias 24/7: ~R\$ 2.842
- **Total infra: R\$ 2.892**

Operacional pós-piloto (se aprovar)

- A100 24/7: R\$ 4.500-5.000
- useapi.net (mantém): R\$ 75
- fal.ai (reduzido, só Hollywood premium): R\$ 300-1.000
- RunPod + Supabase + Hetzner: R\$ 425
- **Total: R\$ 5.300 - 6.500/mês**

Economia projetada

Com volume médio (5k face swaps + 10k FLUX/mês):

- Economia bruta: R\$ 6.000/mês
- Custo A100: -R\$ 2.842/mês
- **Economia líquida: R\$ 3.158/mês**

Payback: 30 dias após piloto.

Com 1000 users ativos: economia de R\$ 6.000/mês líquida, margem operacional sobe de 60% → 85%.

Cronograma (4 semanas)

Semana	Foco	Entrega
1	Provisionar A100 + Docker + ComfyUI base + Cloudflare Tunnel	GET /system_stats retorna info GPU
2	Proxy Node.js + integração <code>provider: "comfyui"</code> + smoke test FLUX schnell	Workflow teste end-to-end via UI
3	Pipeline Face Swap completo (InsightFace + restoration + upscale) + A/B test 10%	Pipeline gera output em <4min, qualidade \geq pixverse
4	FLUX dev workflow + rollout 50% + métricas + documento decisão	GO/NO-GO documentado

Gates GO/NO-GO objetivos

Decisão de expandir baseada em métricas mensuráveis, não opinião:

Qualidade (obrigatório)

- Face Swap \geq pixverse fal.ai em avaliação cega de 20 amostras (votar sem saber qual é qual): self-host vence \geq 10/20
- FLUX dev \geq fal.ai/flux/dev em 20 prompts: self-host vence ou empata \geq 15/20

Estabilidade (obrigatório)

- Uptime A100 \geq 99%
- Taxa de erro pipeline \leq 5%
- Latência p95 Face Swap \leq 5min
- Latência p95 FLUX dev \leq 60s

Financeiro (obrigatório)

- Economia mensal real projetada \geq R\$ 1.500
- Margem operacional do hub não cai

Decisão

Resultado	Ação
3 gates atingidos	GO: expande pra Lipsync + Voice + Motion Control + Wan 2.2 vídeo (8 semanas adicionais)
Qualidade ou Estabilidade falha	ABORT: rollback (feature flag \rightarrow false, A100 deletada). R\$ 2.842 gastos como aprendizado.
Só Financeiro falha	STOP: mantém A100 só pros 2 workflows migrados, não expande

Plano de rollback

- Soft (1h):** feature flag `enable_self_host_*` \rightarrow false. Tráfego volta 100% pra fal.ai. Código mantido.
- Hard (1 dia):** remove templates `*-self` do DB, mata A100, remove case `"comfyui"` do workflow-engine.

- **Custo:** R\$ 0 software + R\$ 2.842 já pago do mês (não recuperável).
-

Por que piloto em vez de migração full 12 semanas

- Valida hipóteses sem comprometer 3 meses do roadmap
 - Investimento upfront 5x menor (R\$ 5k vs R\$ 25k+)
 - Métricas reais em 30 dias justificam ou matam expansão
 - Dev aprende operar A100 + ComfyUI com 10% do tráfego, não 70%
 - Hedge contra mudança rápida do cenário tech AI
-

Perguntas pendentes pro dev

1. Disponibilidade pra 4 semanas focadas (~40h/semana)?
 2. Familiaridade com Docker + GPU containers (curva ~3 dias se nunca usou)?
 3. Já mexeu em ComfyUI workflows complexos (não só consumiu API)?
 4. Comodidade com SSH + ops + monitoring A100 uptime?
 5. Aceita risco de R\$ 2.842 virarem aprendizado se piloto falhar?
 6. Cobertura noite/fim de semana pra rollback rápido se prod quebrar?
-

Decisão pedida

Responder em até 1 semana:

1. **GO / NO-GO no piloto 4 semanas?**
 2. Se GO: data de início e horas/semana dedicadas
 3. Se NO-GO: razão principal (custo / complexidade / prioridade / skill gap) + contraproposta
 4. Riscos adicionais não documentados
 5. Dúvidas técnicas a esclarecer antes de começar
-

Documento técnico completo (774 linhas, com arquitetura detalhada, riscos, cronograma semana-a-semana, specs de hardware e stack): ver [docs/proposals/2026-05-20_hibrido_a100_self_host.md](#) no repo [ghostzuka/knight-s-forge](#) (PR #56).