

# Step Report â€” Stage 0.5: Model Specification

**Data:** 2026-02-17 19:30 BRT **Executor:** Tiuto (Opus 4.6, main session)  
**DuraÃ§Ã£o:** ~15 min **Commit:** 99a5d36

## Objetivo

LLM atua como economista: IÃ³a governance docs + dados histÃ³ricos do OpenClaw e propÃ³ue features repo-especÃ¡ficas para o modelo logit de P(merge).

## Fontes analisadas

Fonte	Tamanho	ConteÃºdo
data/AGENTS.md	21KB	Estrutura do projeto, CI, coding conventions, multi-agency safety, extensÃµes
data/CONTRIBUTING.md	5KB	Lista de 10 maintainers, workflow de contribuiÃ§Ã£o, tÃ©cnica de AI PRs
all_historical_prs.json	3233 PRs	Metadata: labels, author, additions/deletions, draft, milestone
enriched_full.jsonl	1521 PRs	Comments, reviews, files detalhados
D01v3-enrichment-analysis.md	Report	Achados estatÃsticos do Stage 0

## Artefatos entregues

Artefato	Path	ConteÃºdo
Model Spec	model_spec.json (v0.5.0)	33 features, split early/mature, estratÃ©gia temporal dual
Spec Notes	docs/model-spec-notes.md	Racional: insights dos governance docs, exclusÃµes com motivo
Feature Map	features/feature_map.json	Mapeamento operacional: feature â†’ regra de extraÃ§Ã£o
Key-Info	key-info.json (K010-K013)	4 entries novas

## Features especificadas

### Early Model (19 features â€” disponÃ¡veis na criaÃ§Ã£o da PR)

Feature	Tipo	Sinal esperado	Fonte
loc_additions	numÃ©rica		PR metadata

Feature	Tipo	Sinal esperado	Fonte
		negativo após threshold	
loc_deletions	numérica	positivo/ neutro	PR metadata
loc_total	numérica	negativo após threshold	computado
files_changed	numérica	negativo	PR metadata
size_label	categórica	não-linear (S/M/L > XS/ XL)	labels
has_tests	binária	positivo	file paths (*.test.ts)
ci_green	binária	positivo	áσ i, ♦ GitHub checks API (gap)
is_draft	binária	negativo	PR metadata
category	categórica	varia por tipo	labels (agents/ gateway/ cli/docs...)
component_area	categórica	varia	labels (channel:/ app:) + paths
author_prior_prs	numérica	positivo	computado (temporal guard)
author_prior_merge_rate	numérica	positivo	computado (temporal guard)
author_association	categórica	MEMBER > CONTRIBUTOR > NONE	enriched
has_maintainer_label	binária	<b>forte positivo (90.7%)</b>	labels
has_trusted_contributor_label	binária	positivo	labels
has_experienced_contributor_label	binária	positivo	labels

Feature	Tipo	Sinal esperado	Fonte
weeks_since_open	numérica	positivo (confounding)	computado
weekly_pr_volume	numérica	negativo (controle)	computado
release_period	categórica	varia	GitHub tags API

## Mature Model (+14 features â€” interação acumulada)

Feature	Tipo	Sinal esperado	Fonte
comment_count	numérica	não-linear	enriched
high_engagement	binária	positivo (4+ comments)	computado
review_count	numérica	positivo	enriched
has_approval	binária	positivo	enriched (review state)
has_changes_requested	binária	negativo	enriched (review state)
has_maintainer_comment_stipe	binária	<b>forte positivo (74.5%, 17x)</b>	enriched
has_top_contributor_comment	binária	positivo	enriched (authorAssociation)
top_contributor_comment_count	numérica	positivo	enriched
has_greptile_review	binária	fraco/negativo (controle)	enriched
greptile_score	numérica	fraco/zero	enriched
pr_age_hours	numérica	positivo â†' plateau	computado
touches_multiple_channels	binária	negativo	labels + paths
touches_extensions	binária	varia	labels + paths
is_fork_pr	binária	negativo	âš i GitHub API (gap)

## Achados repo-específicos

---

1. **CLAUDE.md** é symlink para **AGENTS.md** → não são docs separados
2. **maintainer\_label = 90.7% merge** (183 PRs) → sinal mais forte por label, disponível no early model
3. **Labels trusted-contributor (96)** e **experienced-contributor (99)** → sinais de status do contribuidor
4. **AGENTS.md adverte explicitamente** sobre refactoring cross-channel → touches\_multiple\_channels
5. **Extensions têm regras de packaging prioritárias** → touches\_extensions como sinal distinto
6. **1318 autores únicos, distribuídos heavy-tail** → history de autor é esparsa para maioria
7. **requested\_reviewers quase nunca usado** (5 PRs) → excluído

## Exclusões (com motivo)

---

Feature	Motivo
requested_reviewers	5 PRs apenas. Zero poder discriminativo
milestone	Não explorado. Deferido pra v2
is_ai_generated	Precisaria NLP no body. Compliance desconhecida. v2
pr_template_compliance	NLP comparativo com template. v2

## Gaps pra resolver no Stage 1

---

1. **ci\_green** → não está no dataset. Precisa GitHub checks API no ingest
2. **is\_fork\_pr** → não está no dataset. Precisa head.repo vs base.repo no ingest

## Viabilidade

---

- **100% das features têm mapeamento executável** (feature\_map.json)
- **0 feature sem origem identificada**
- **2 features precisam de dados adicionais** (ci\_green, is\_fork\_pr) → resolver no Stage 1
- **Cobertura de enrichment: 47%** (1521/3233) → suficiente para estimar

## Próximo passo

---

Stage 0.7 (Bootstrap Sequencial) ou Stage 1 (Ingest PRs abertas) → parallelizáveis.

---

*Report gerado conforme processo definido no PLAN-v4-DIFF.md.*