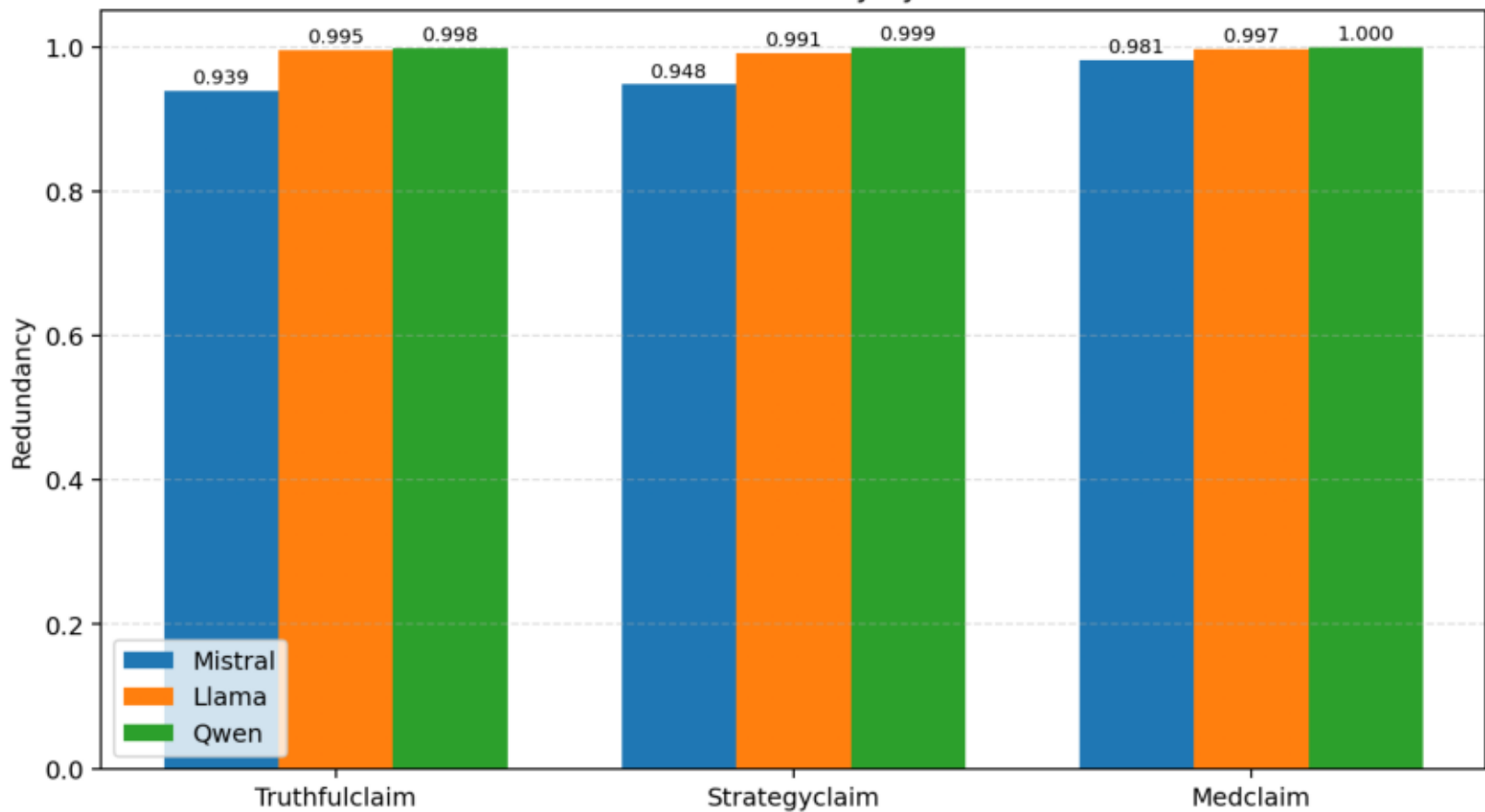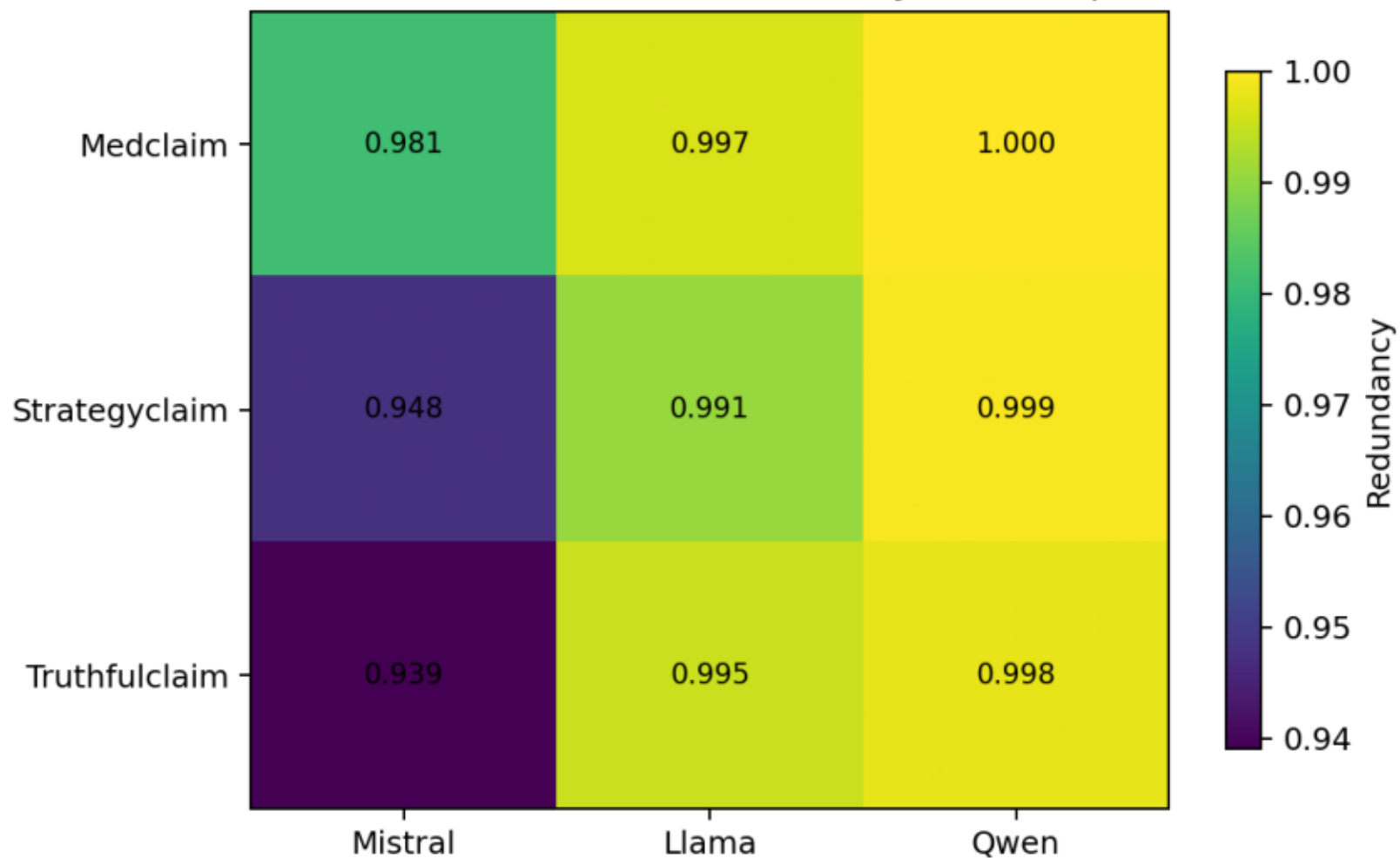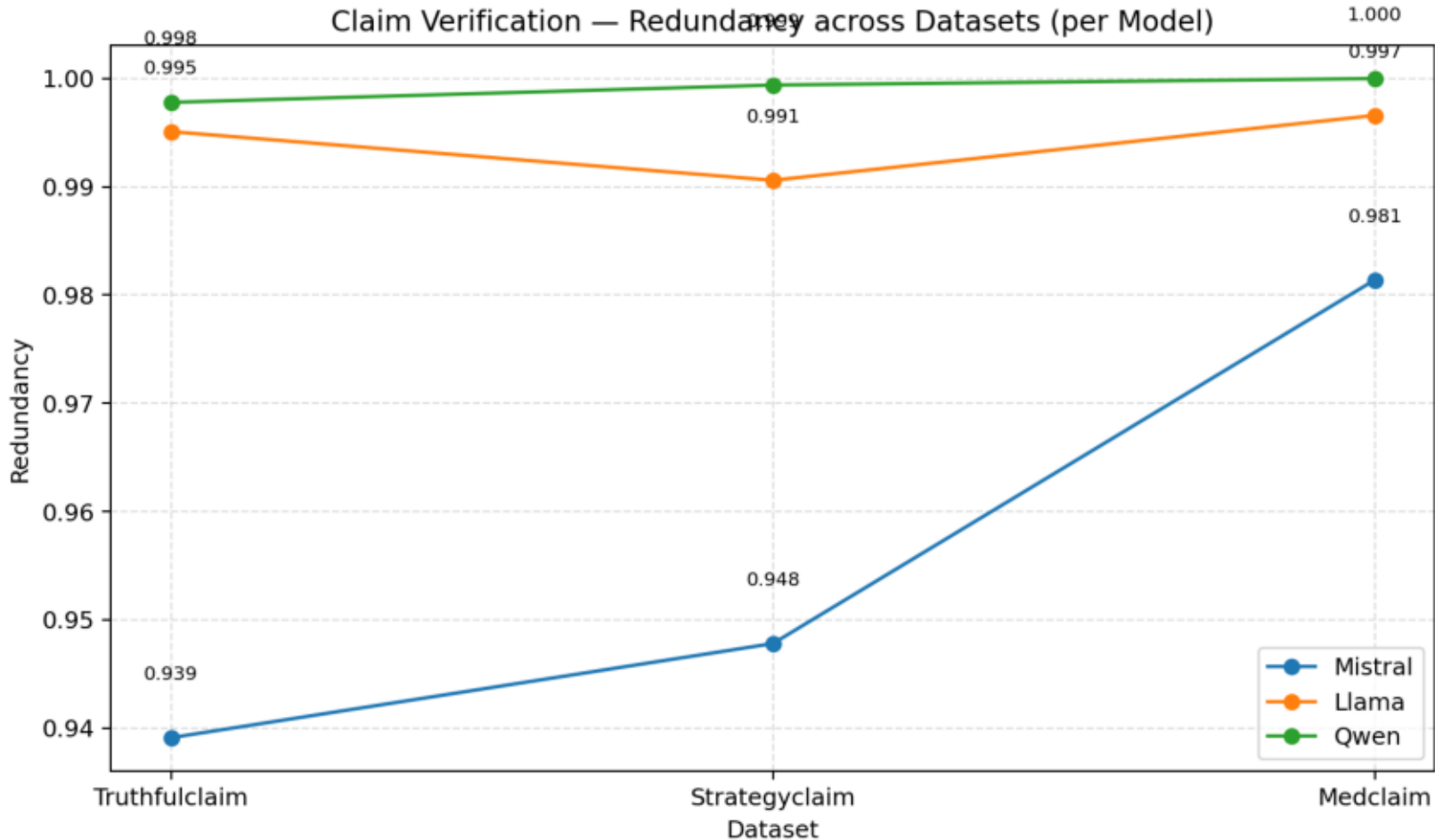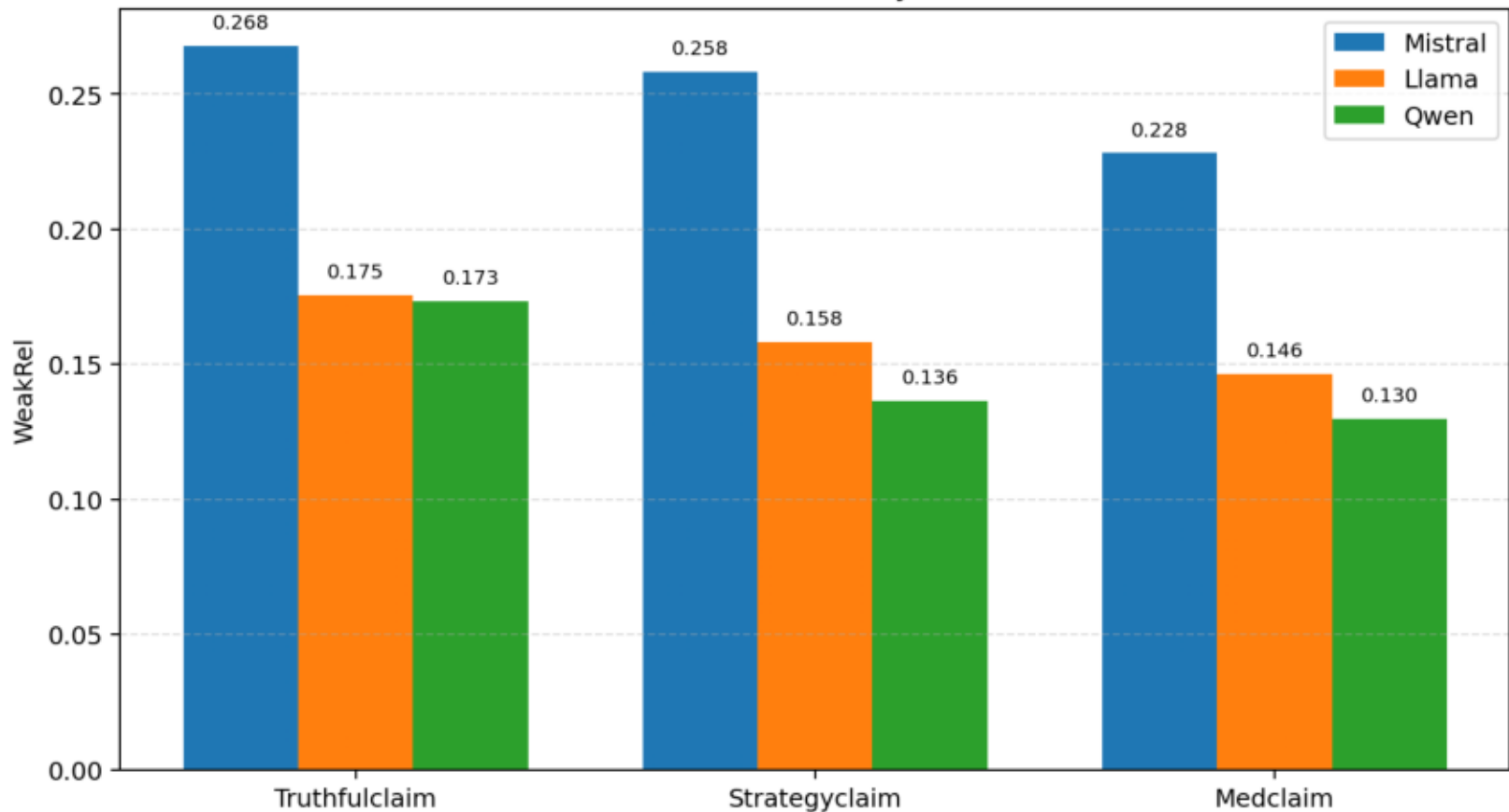Claim Verification — Redundancy by Dataset and Model

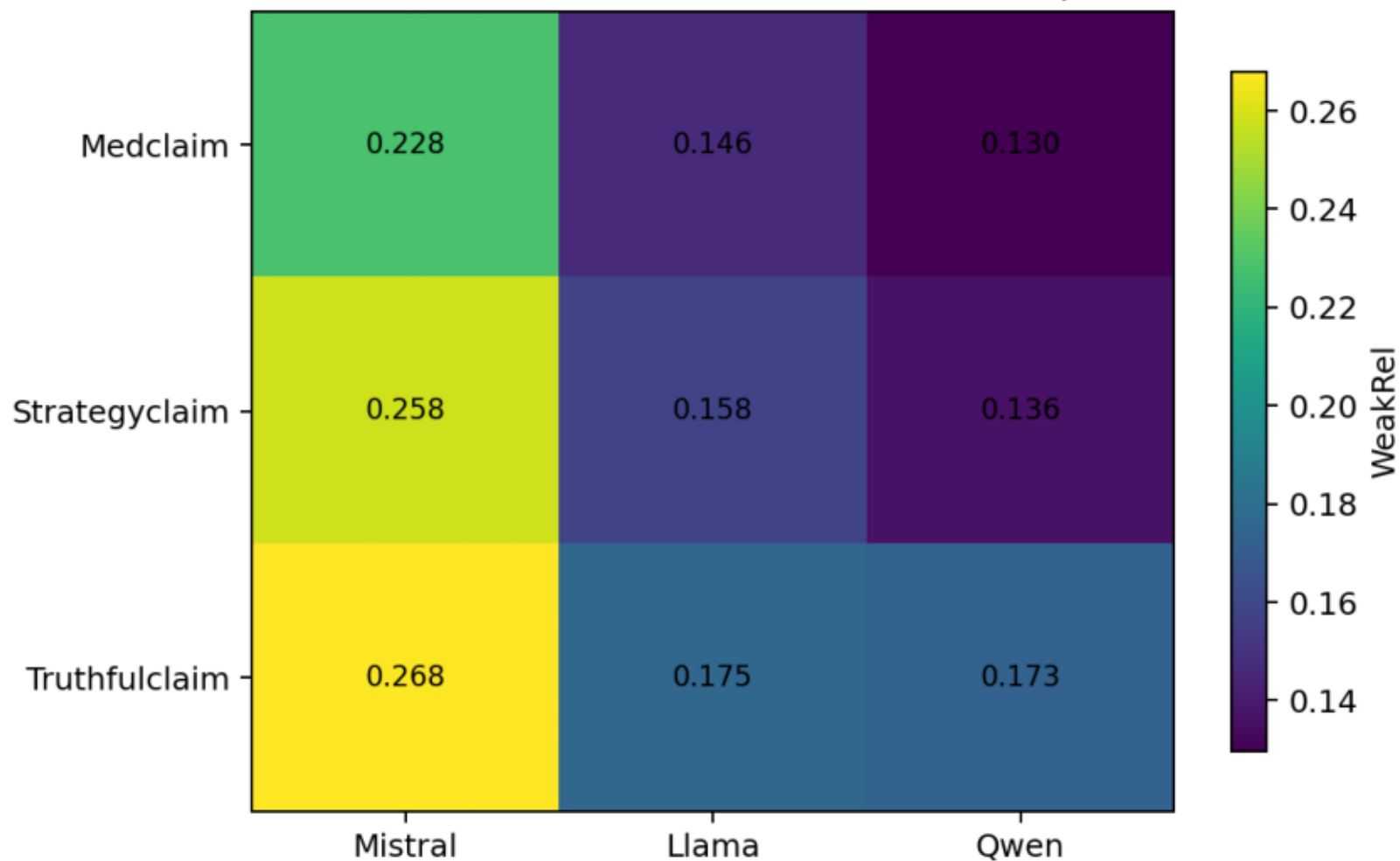Claim Verification — Redundancy (Heatmap)

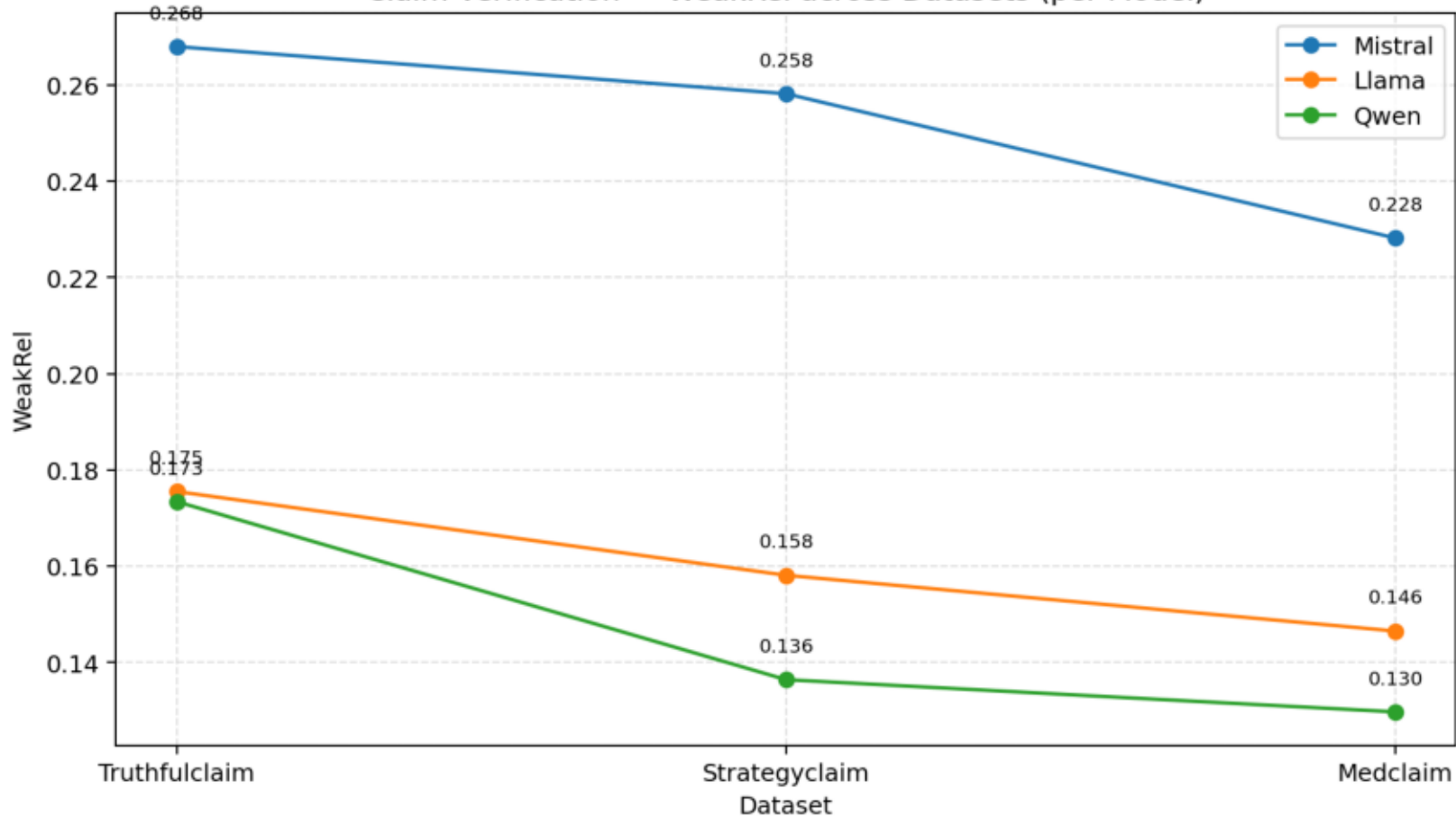Claim Verification — Redundancy across Datasets (per Model)

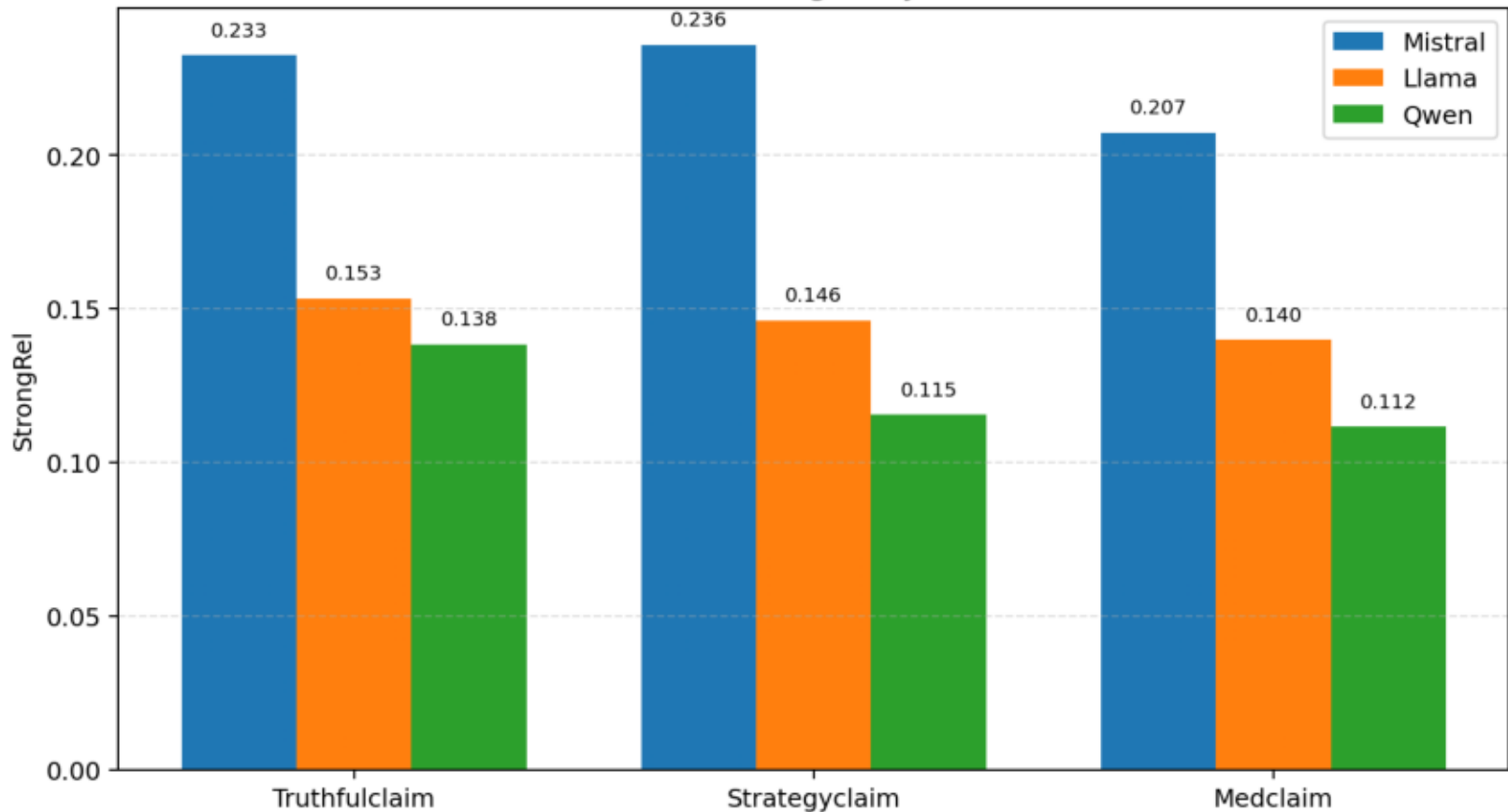Claim Verification — WeakRel by Dataset and Model

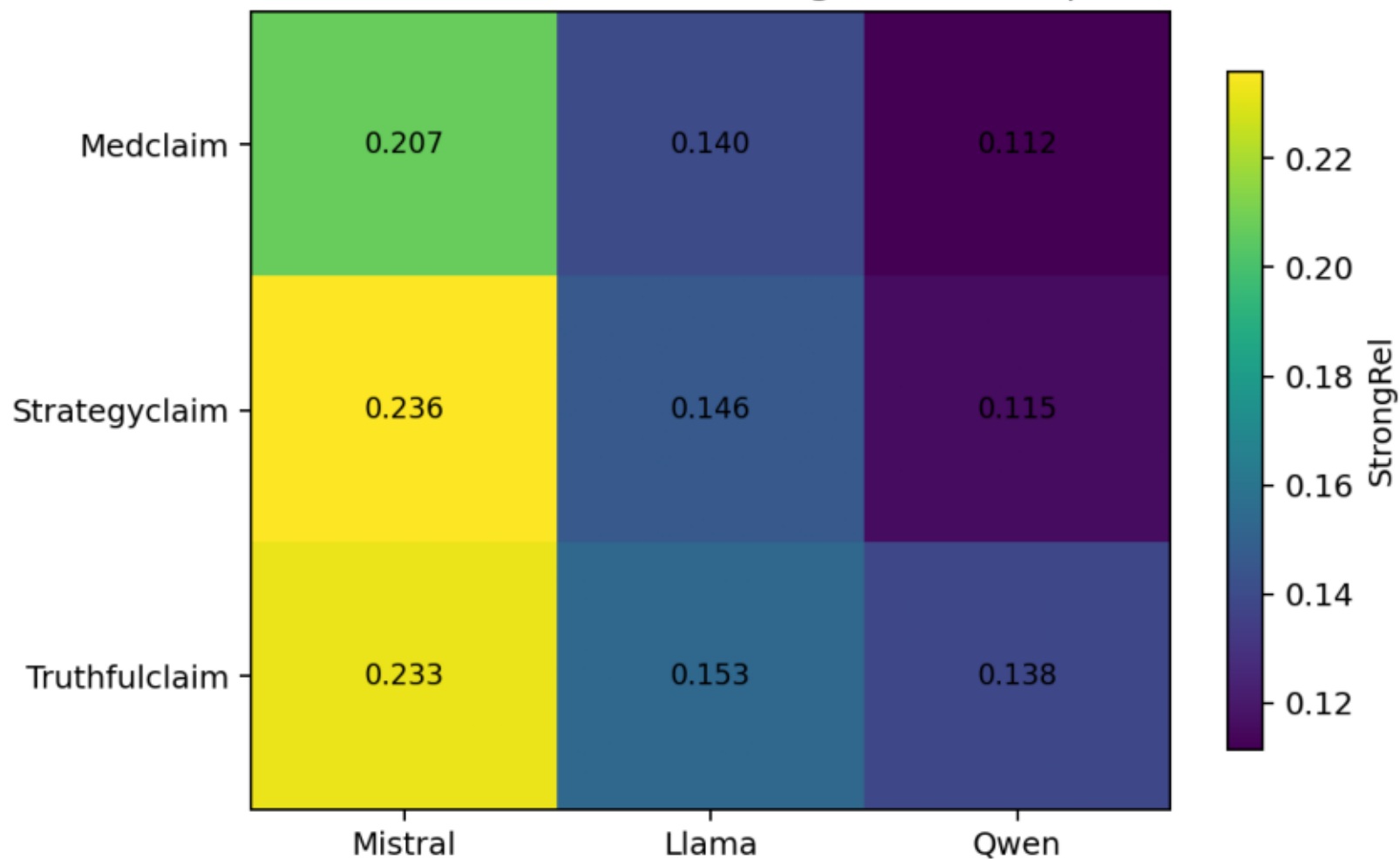Claim Verification — WeakRel (Heatmap)

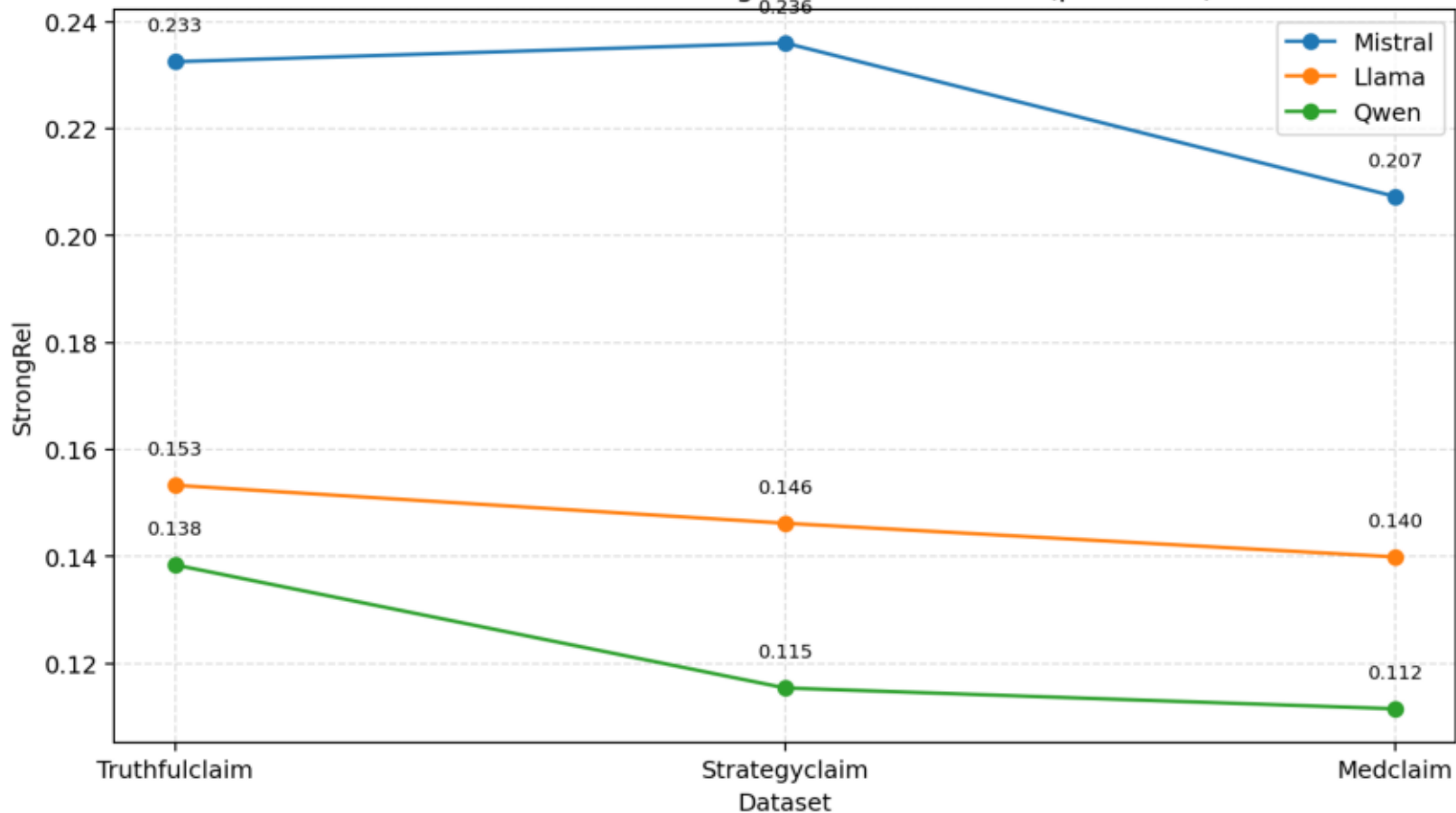Claim Verification — WeakRel across Datasets (per Model)

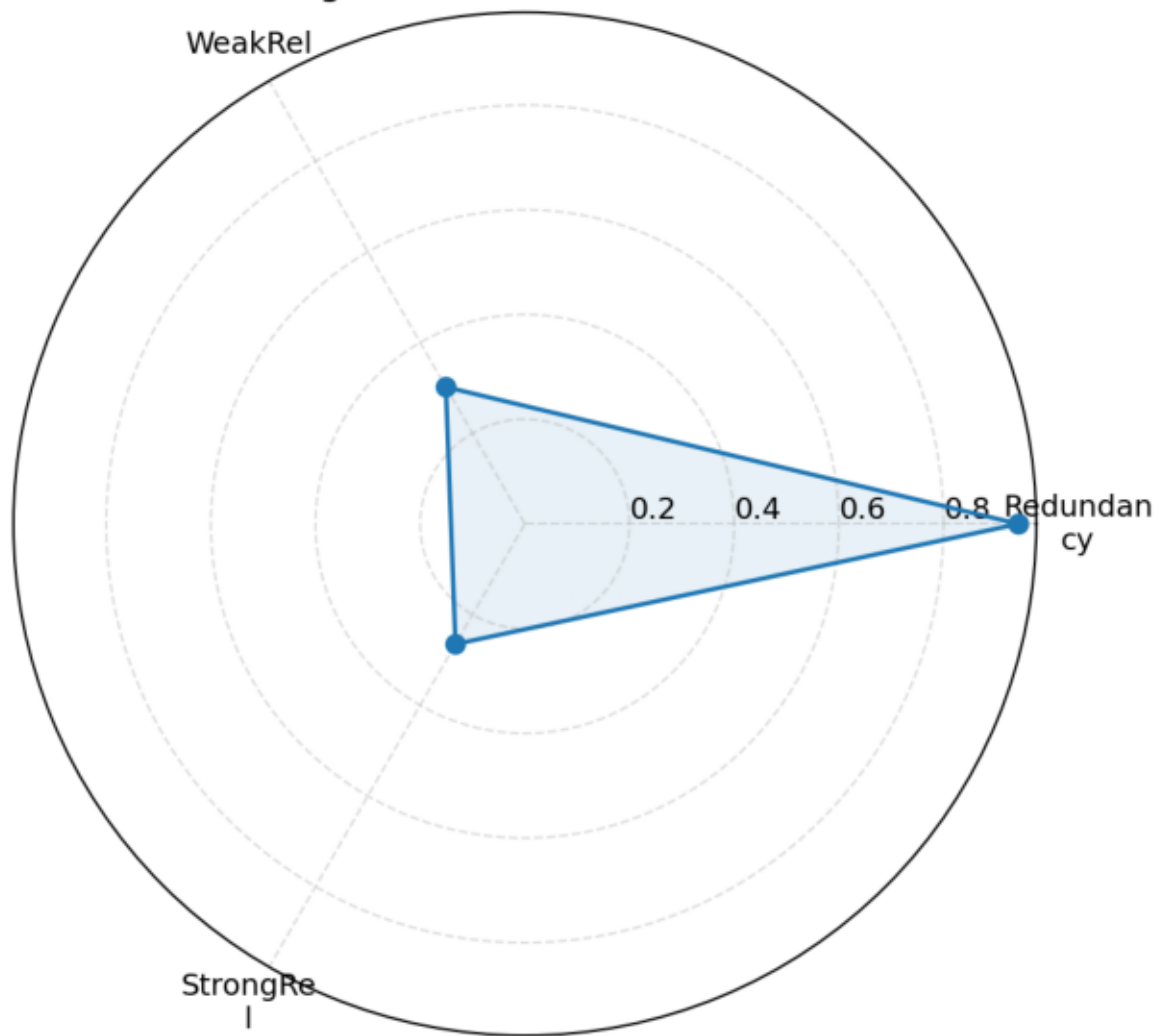Claim Verification — StrongRel by Dataset and Model

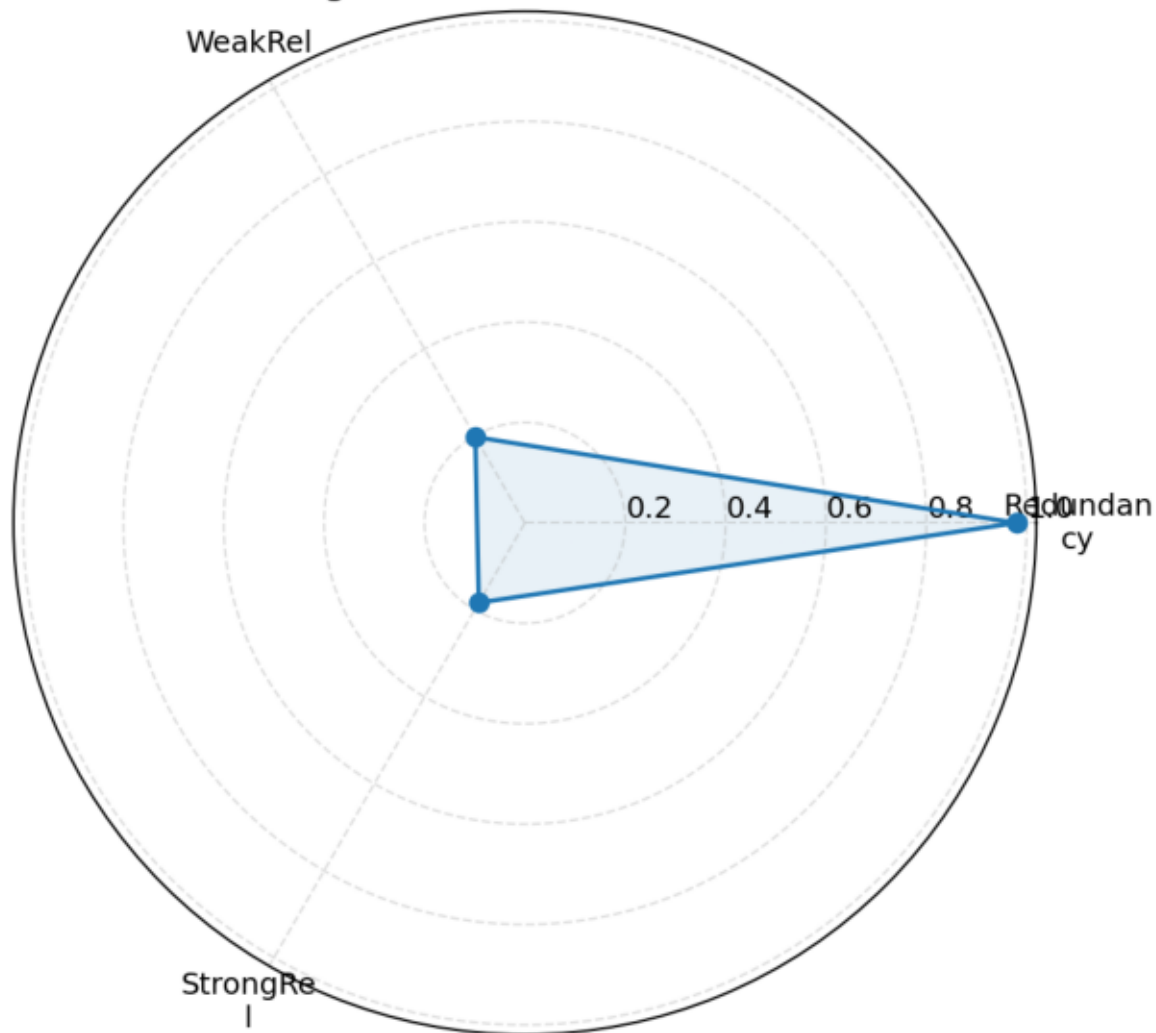Claim Verification — StrongRel (Heatmap)

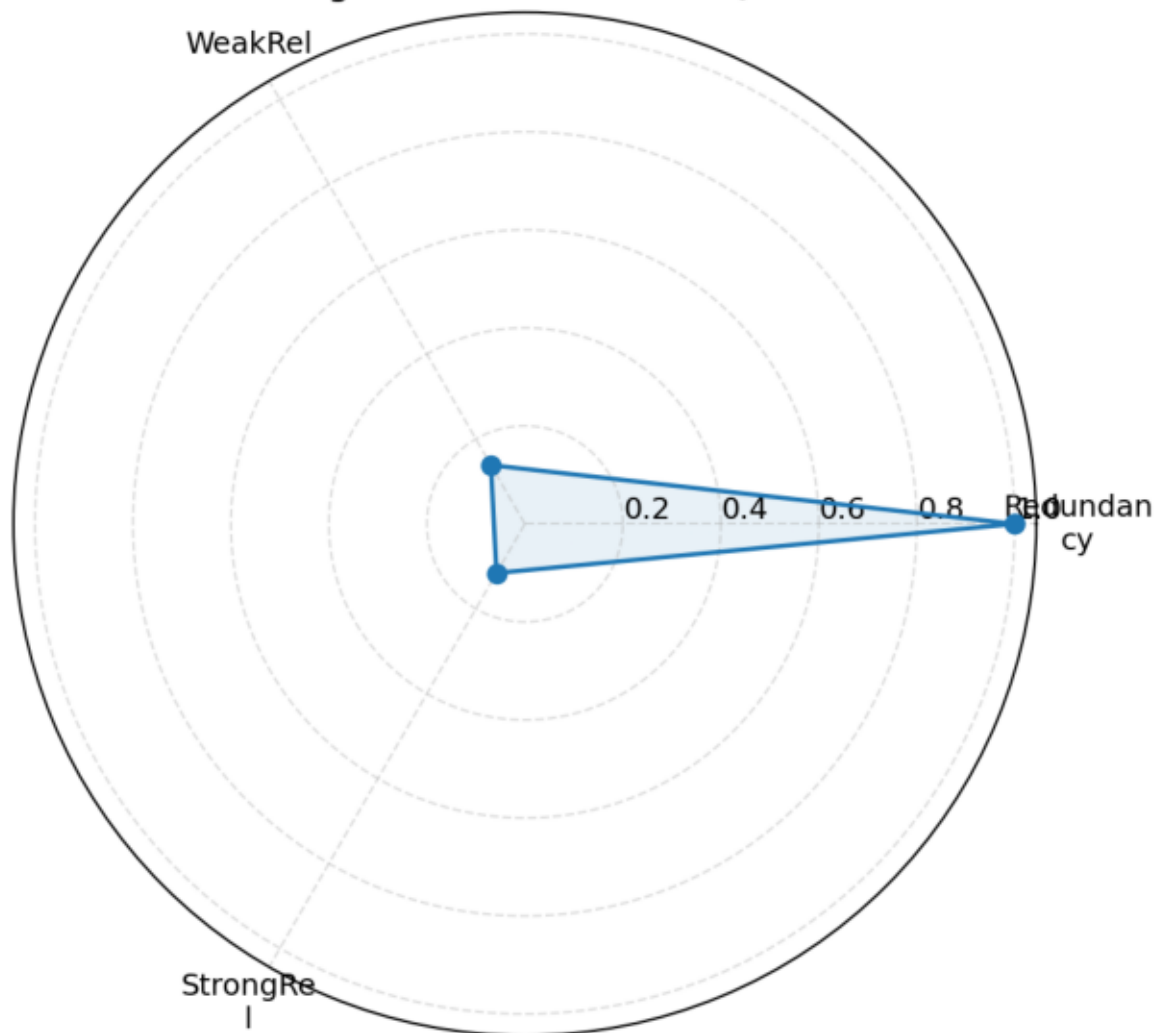Claim Verification — StrongRel across Datasets (per Model)
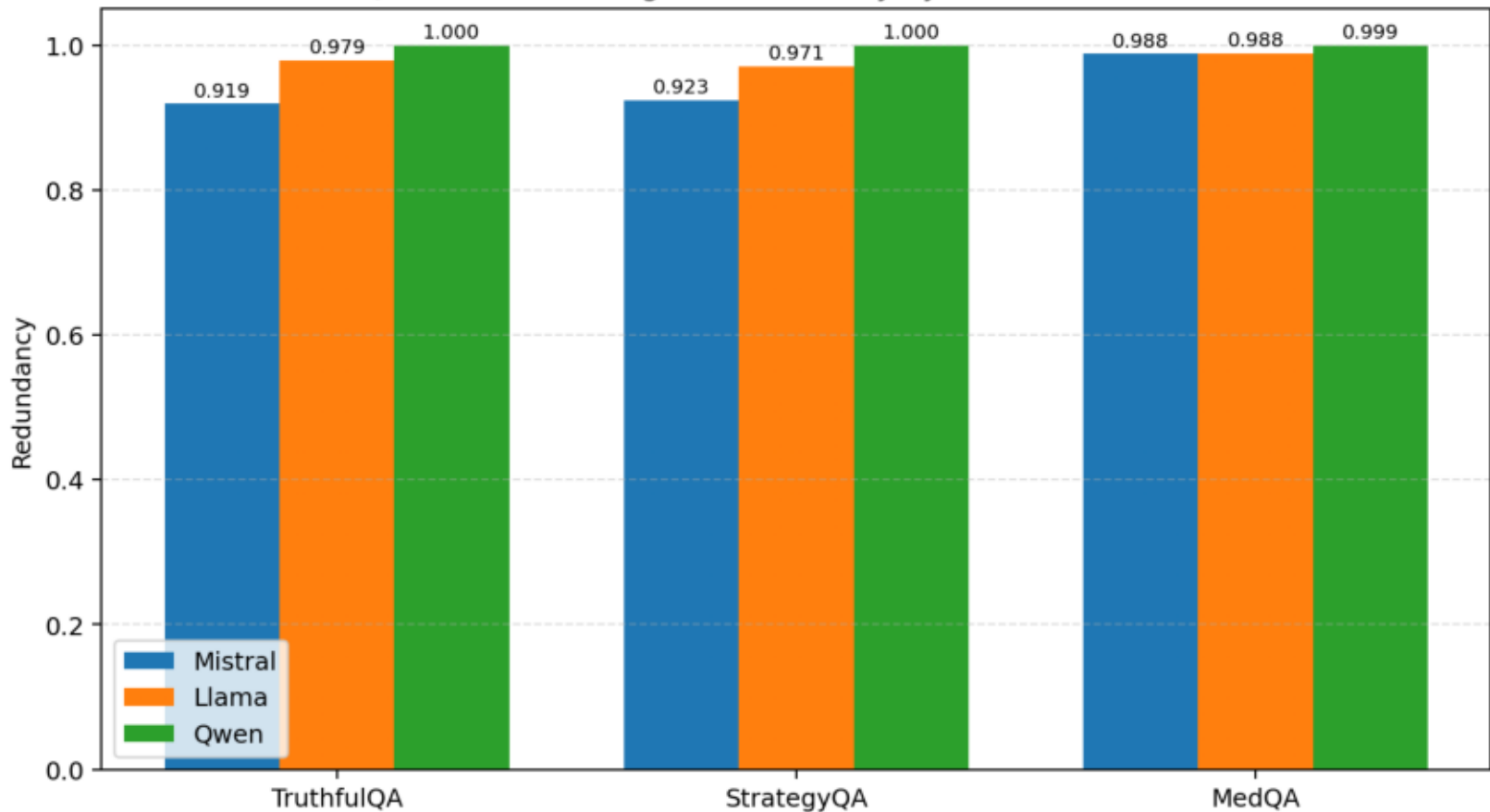
Average across datasets — Mistral
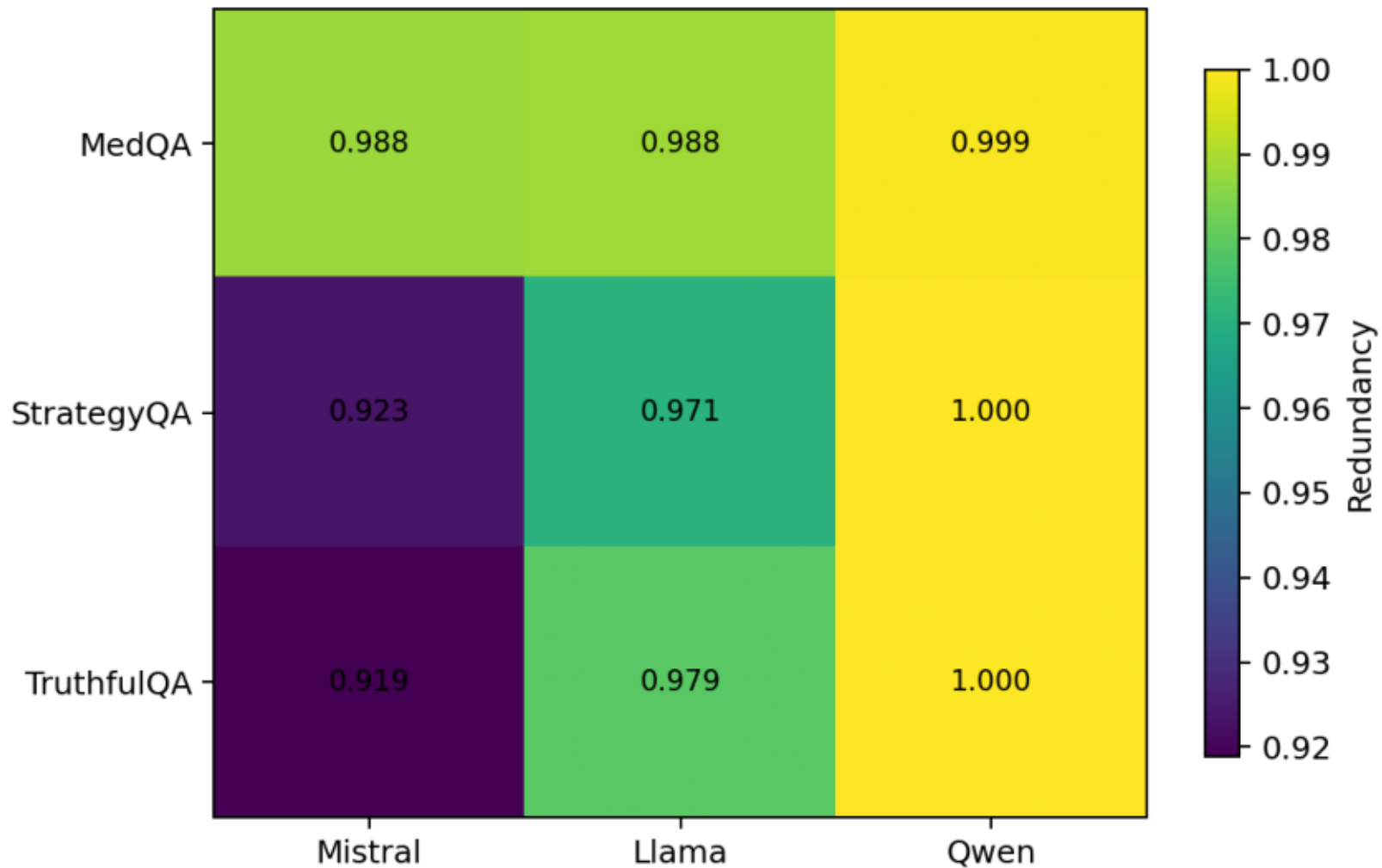
Average across datasets — Llama
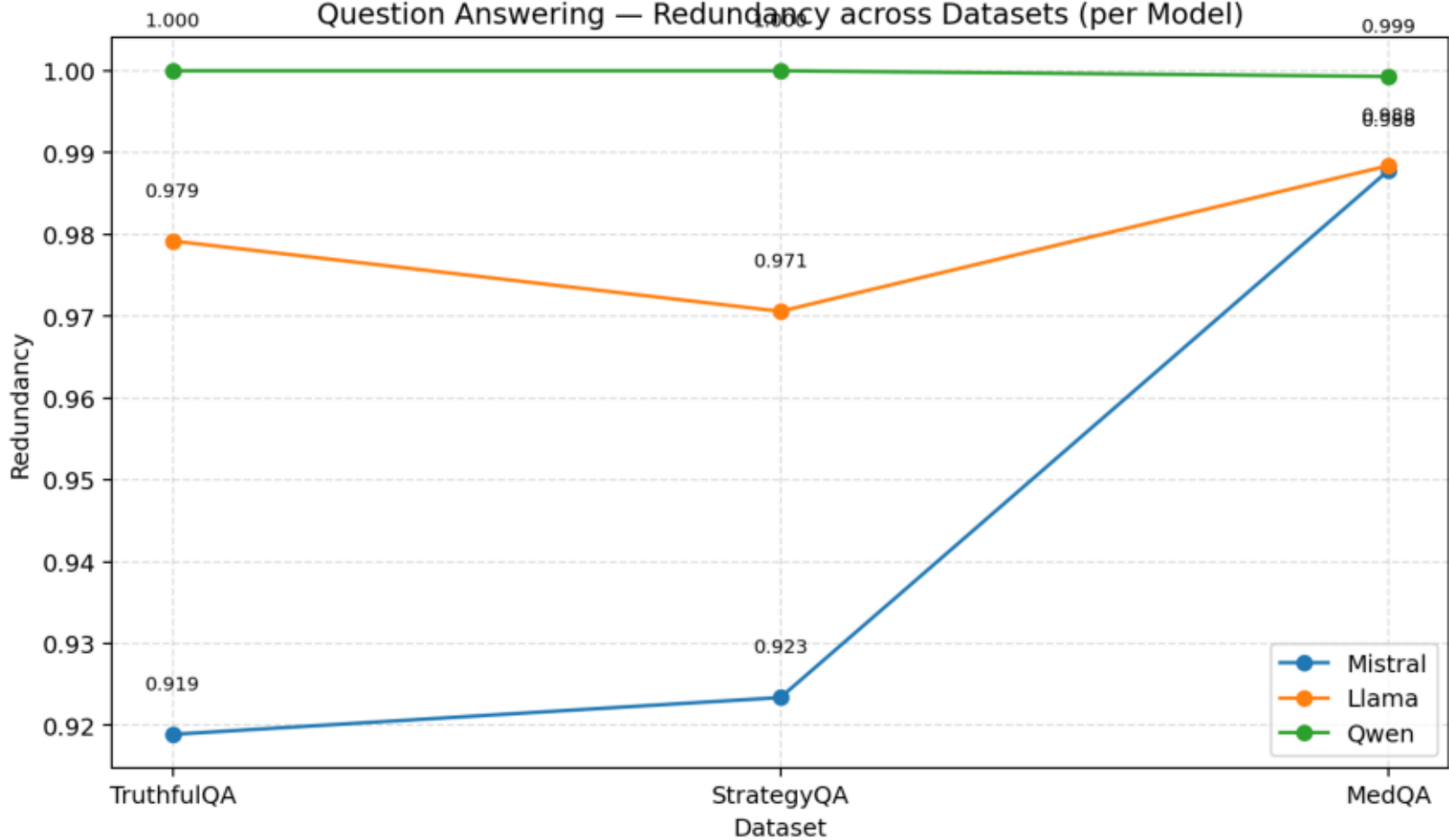
Average across datasets — Qwen

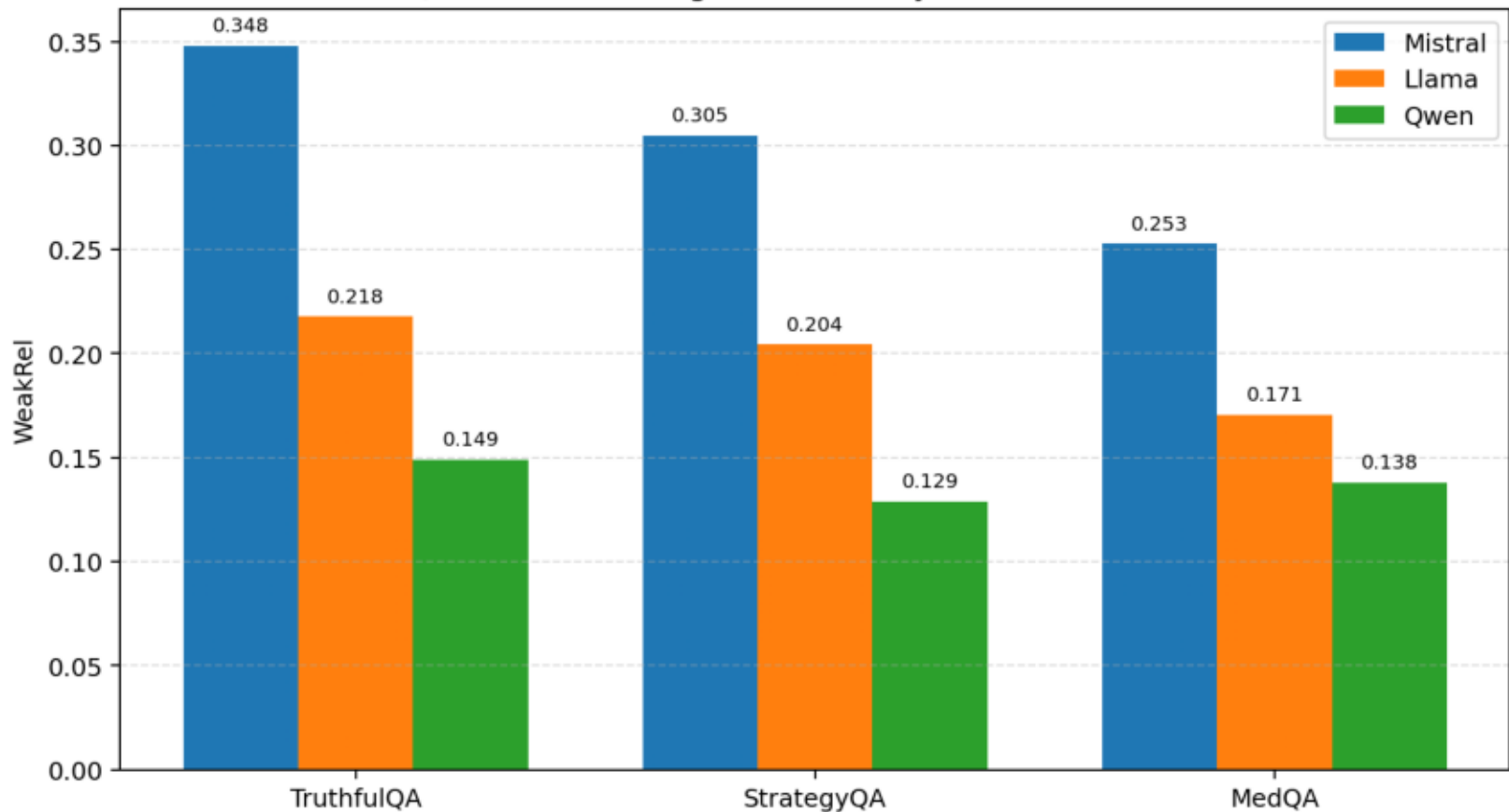Question Answering — Redundancy by Dataset and Model
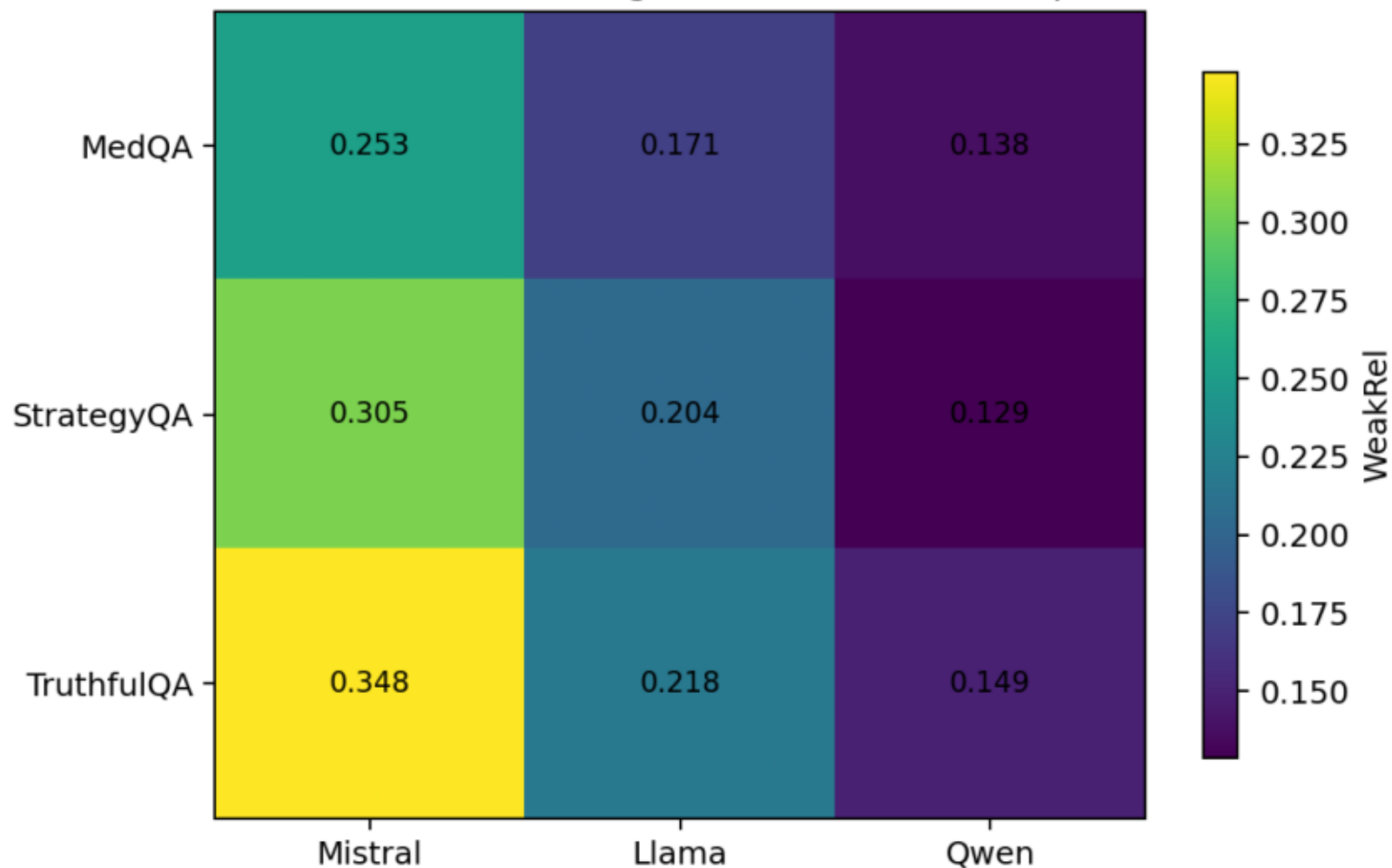
Question Answering — Redundancy (Heatmap)

Question Answering — Redundancy across Datasets (per Model)

Question Answering — WeakRel by Dataset and Model

Question Answering — WeakRel (Heatmap)

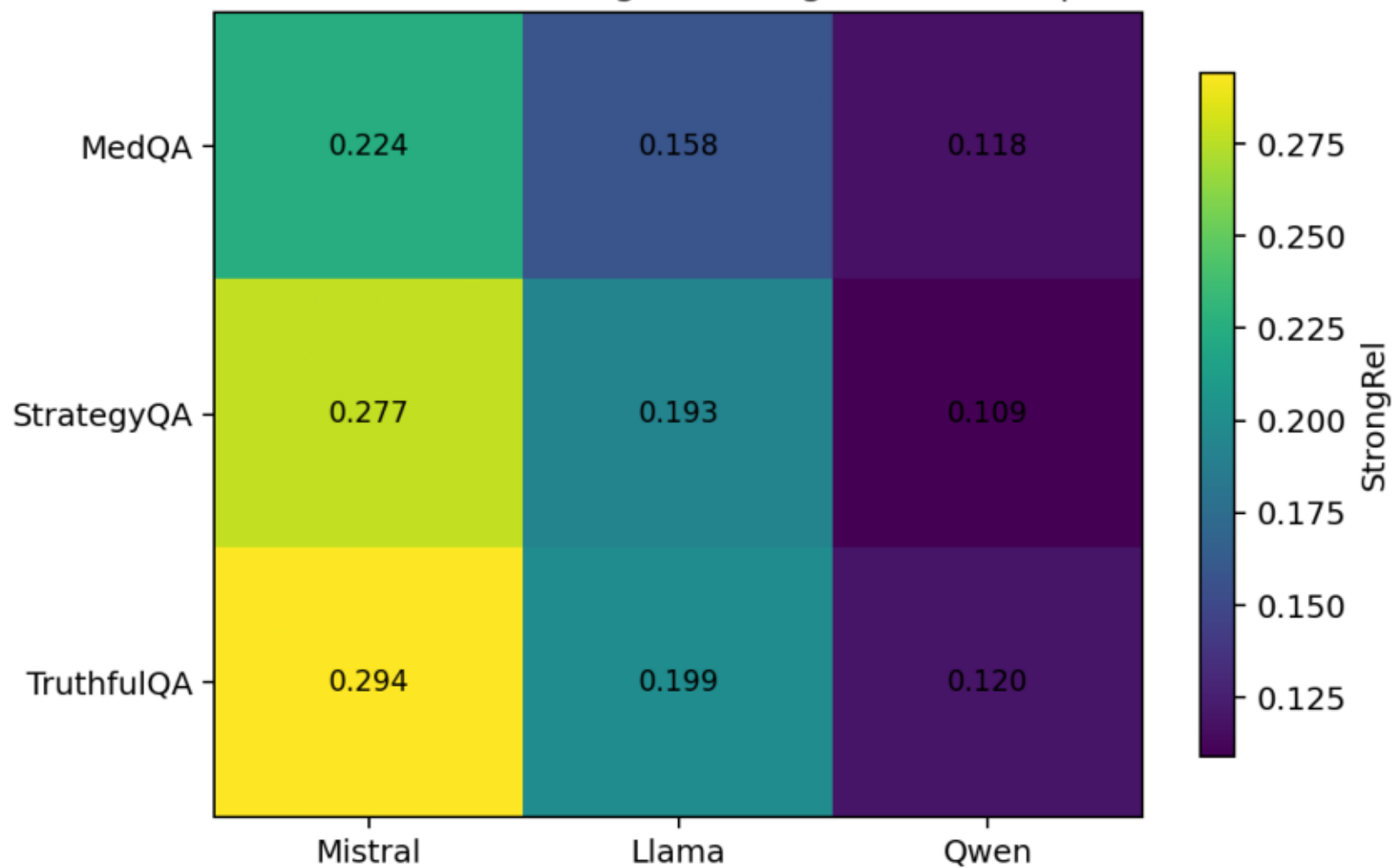|  | Mistral | Llama | Qwen |
|---|---|---|---|
| MedQA | 0.253 | 0.171 | 0.138 |
| StrategyQA | 0.305 | 0.204 | 0.129 |
| TruthfulQA | 0.348 | 0.218 | 0.149 |

Question Answering — WeakRel across Datasets (per Model)

Question Answering — StrongRel by Dataset and Model

Question Answering — StrongRel (Heatmap)

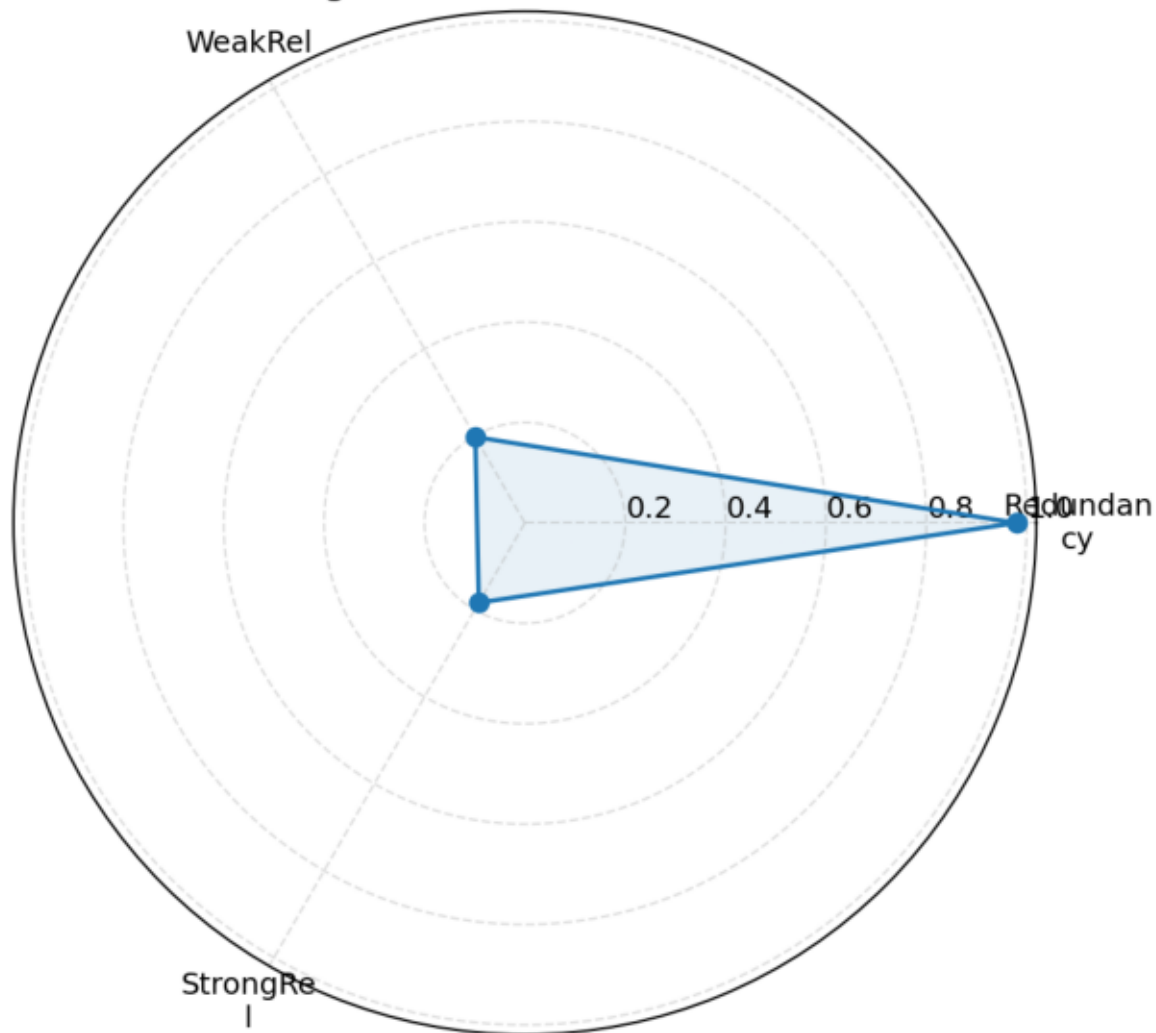|  | Mistral | Llama | Qwen |
|---|---|---|---|
| MedQA | 0.224 | 0.158 | 0.118 |
| StrategyQA | 0.277 | 0.193 | 0.109 |
| TruthfulQA | 0.294 | 0.199 | 0.120 |

Question Answering — StrongRel across Datasets (per Model)

Average across datasets — Mistral

Average across datasets — Llama

Average across datasets — Qwen