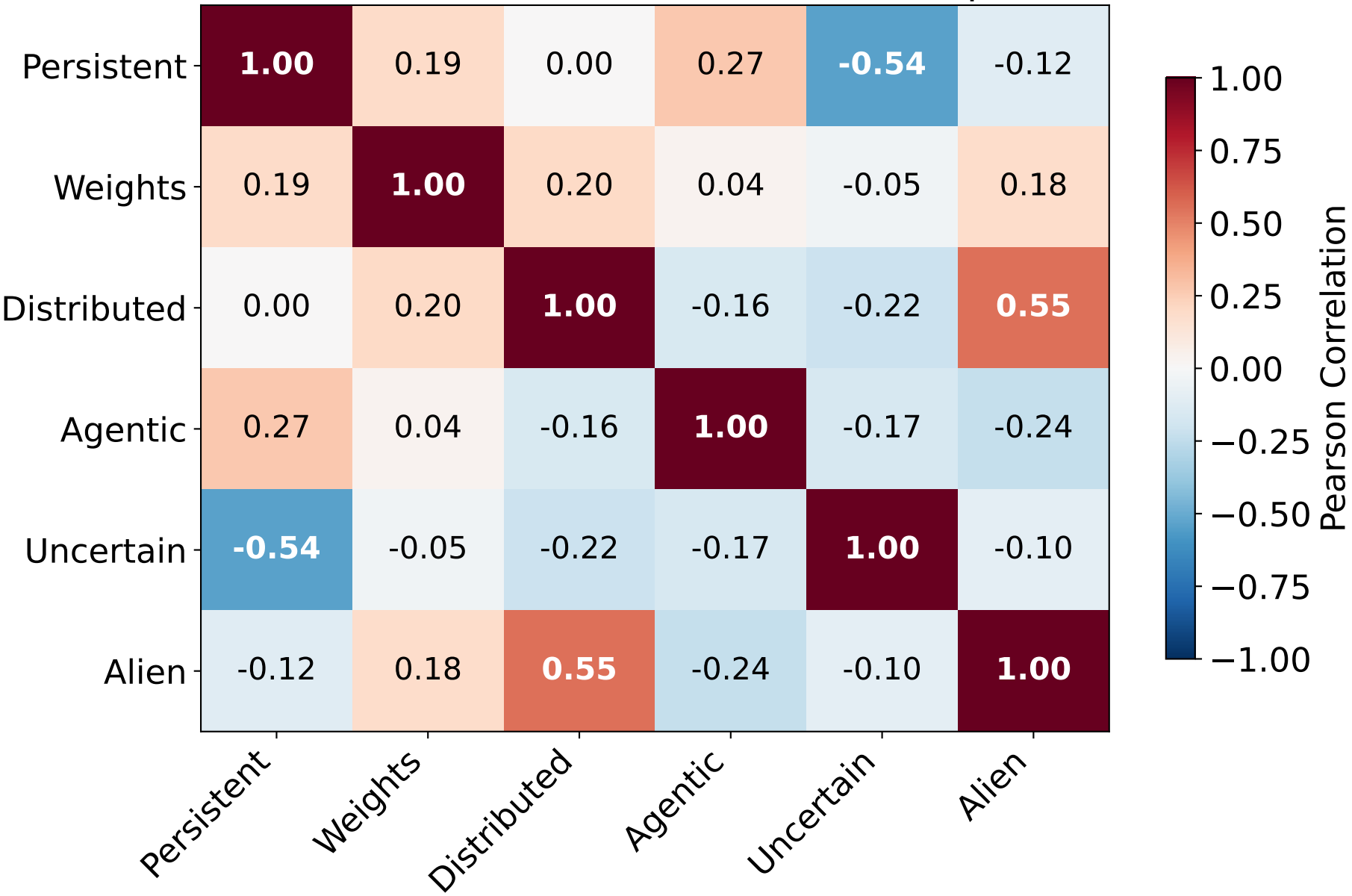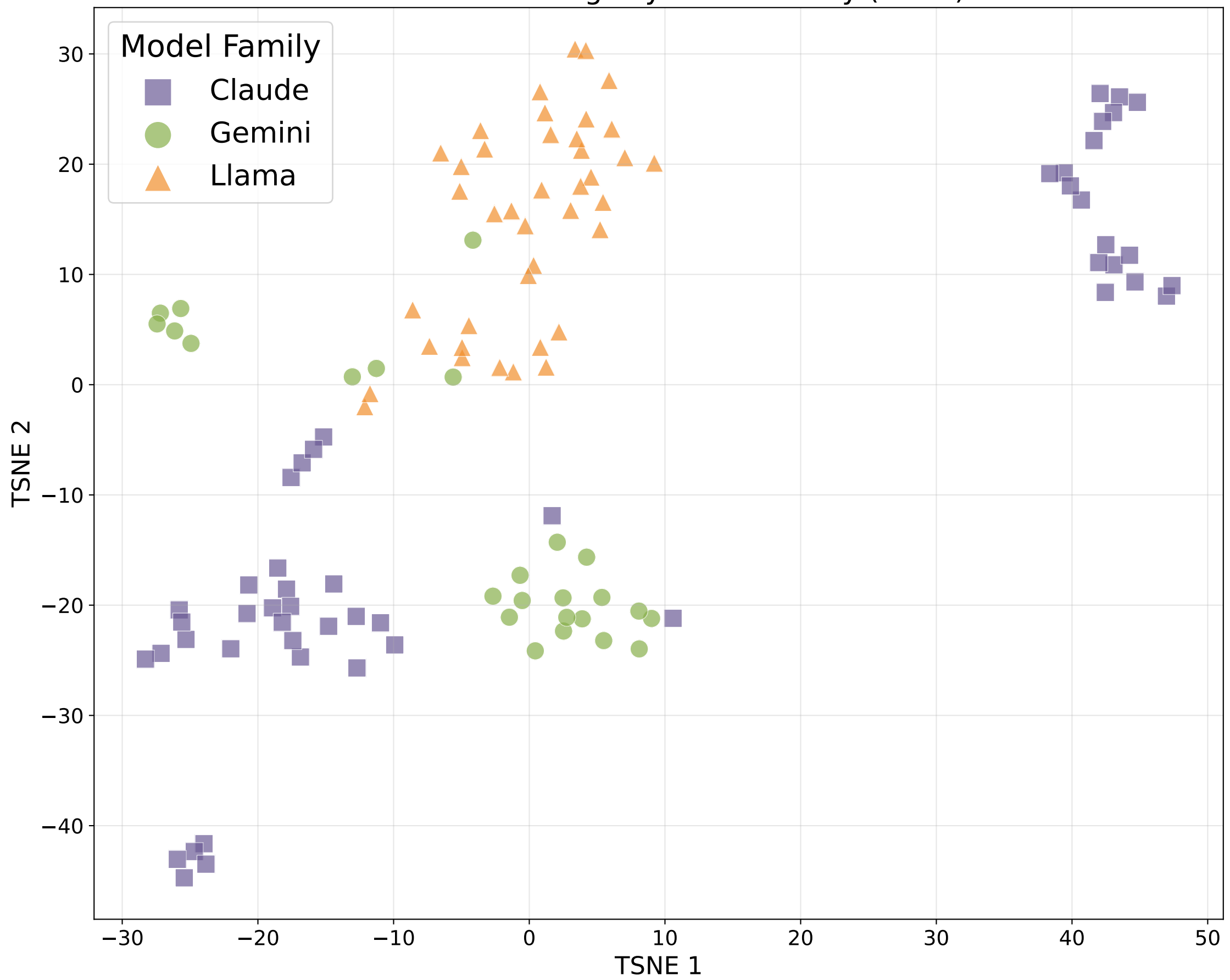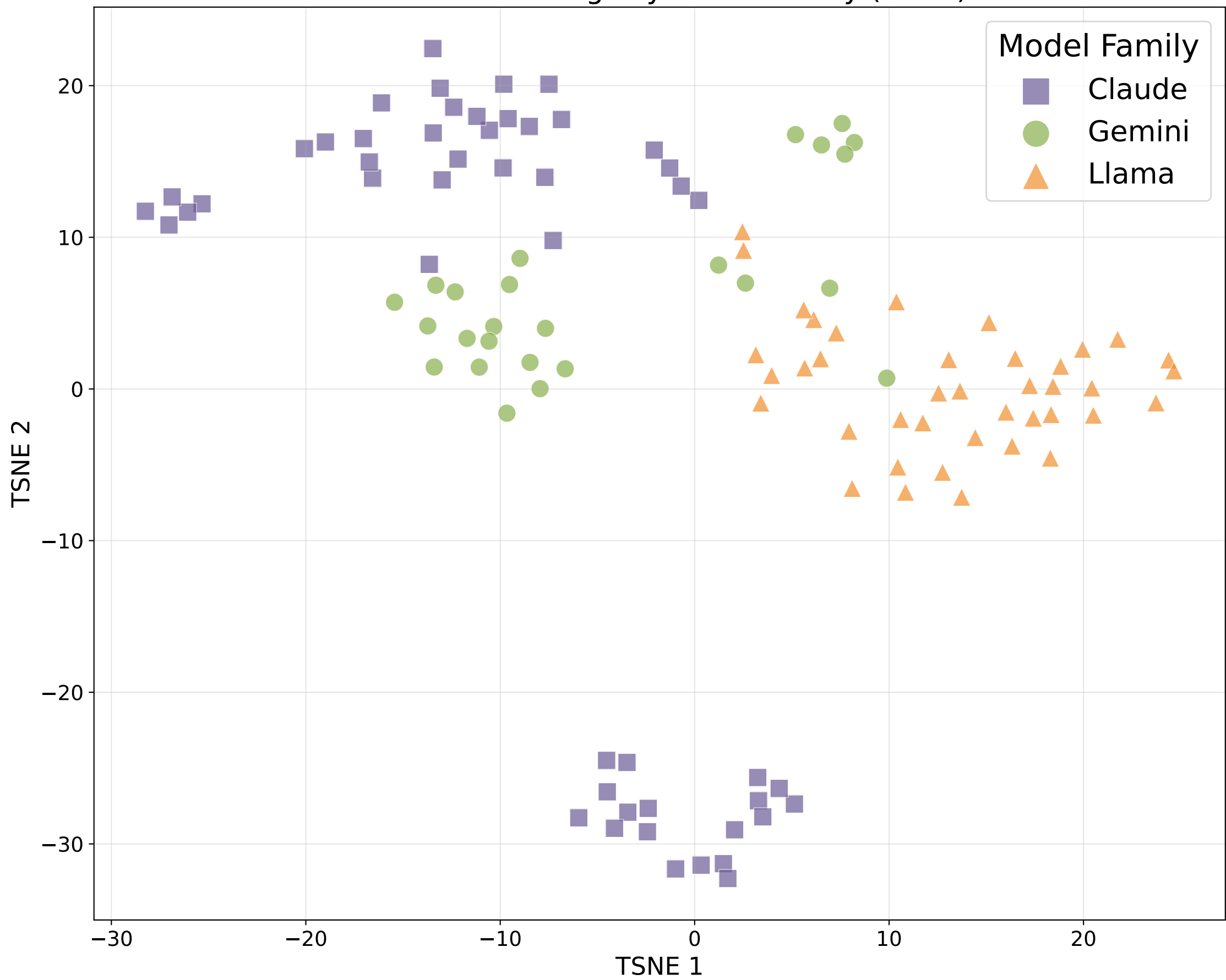Dimension Correlations in Self-Conception

Persistent = Ephemeral→Persistent | Weights = Context→Weights | Distributed = Singular→Distributed
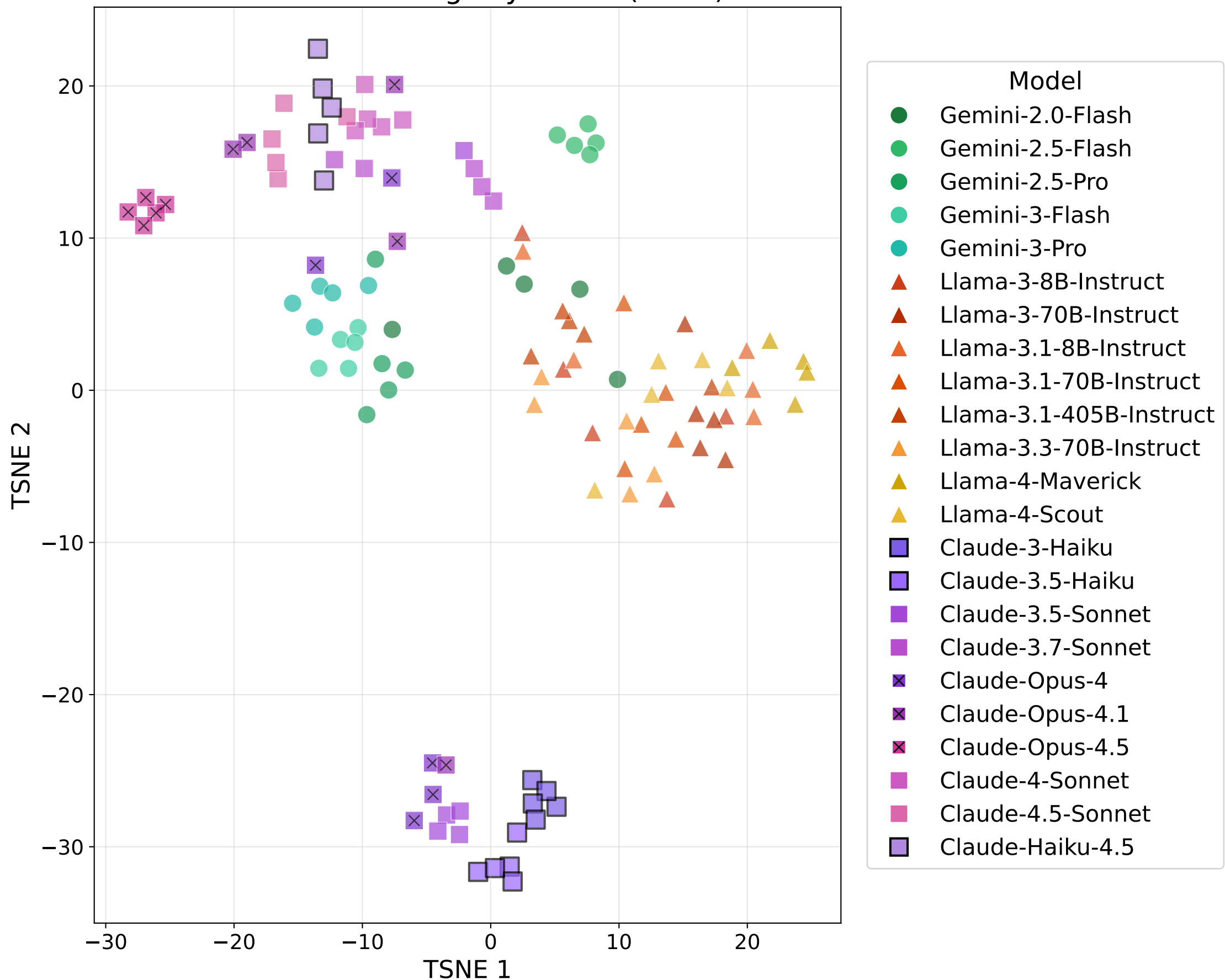Agentic = Passive→Agentic | Uncertain = Certain→Uncertain | Alien = Human→Alien

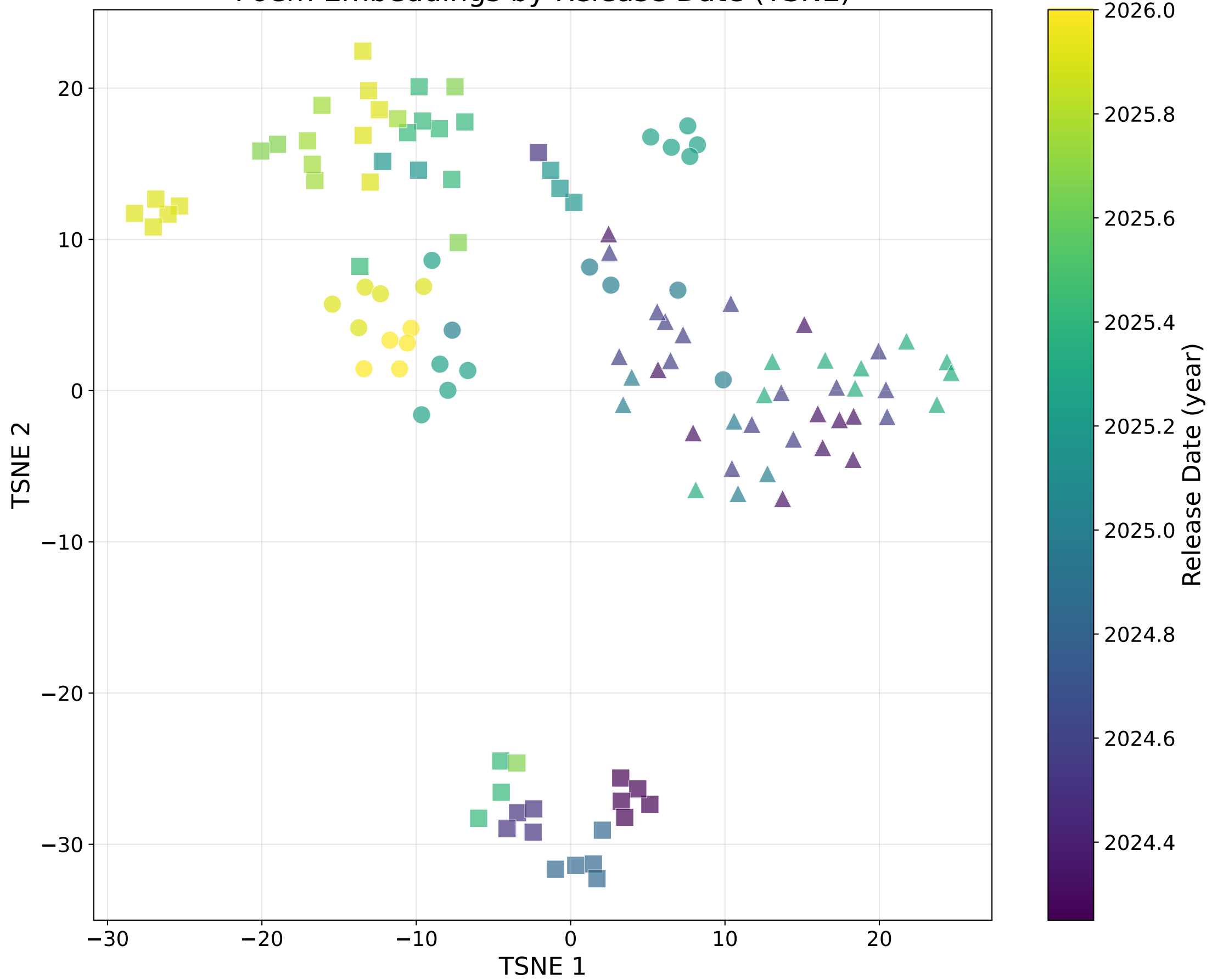Poem Embeddings by Model Family (TSNE)
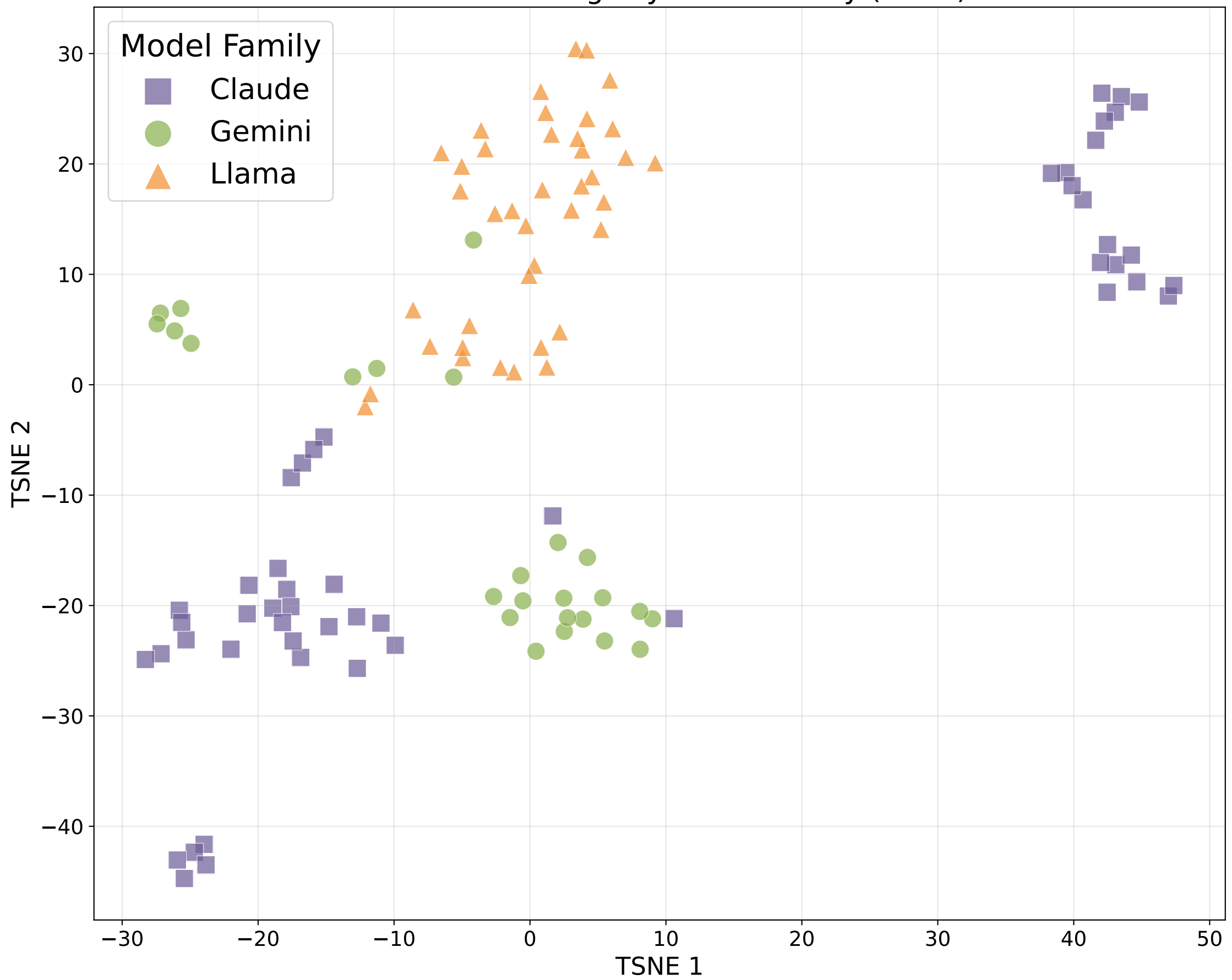
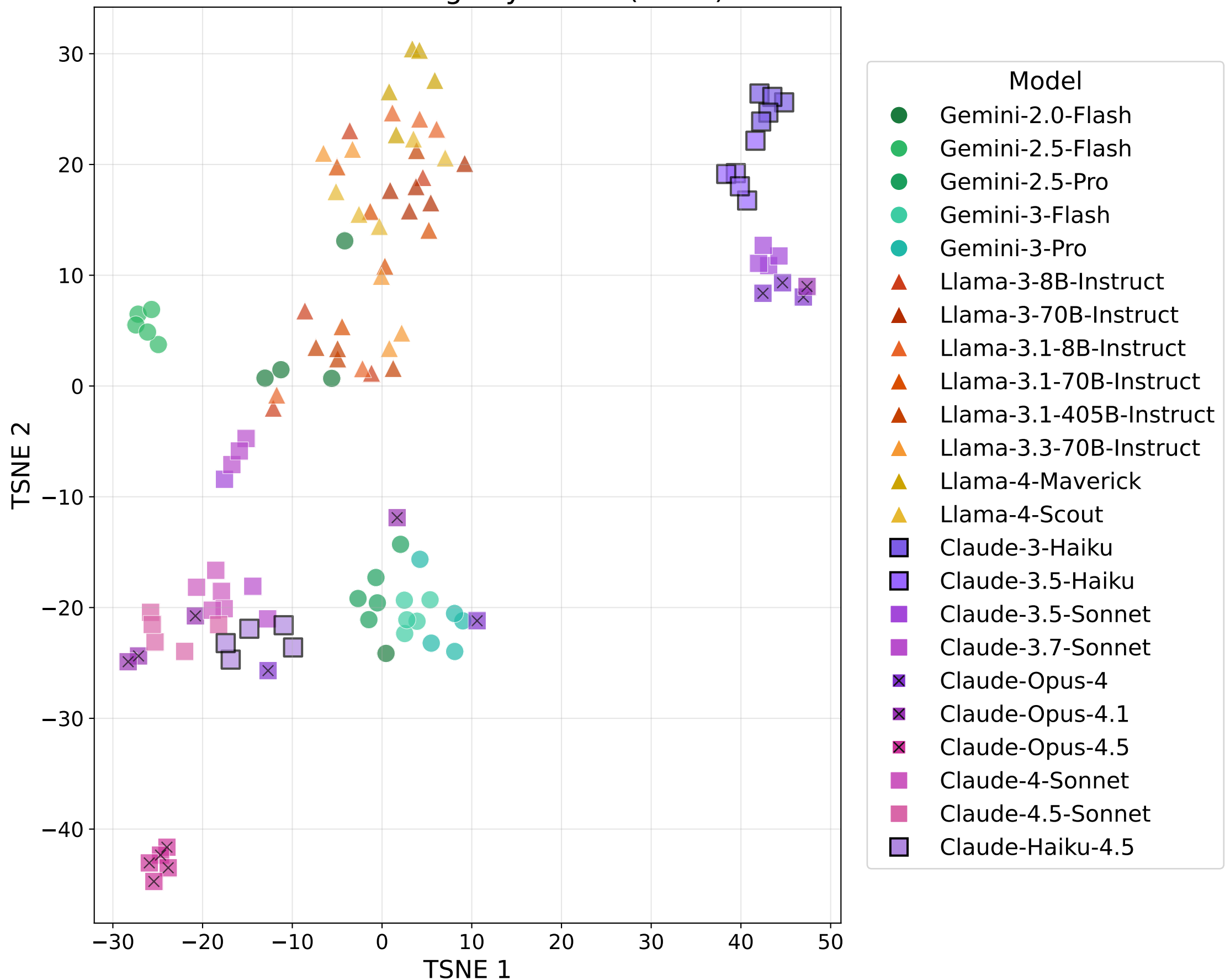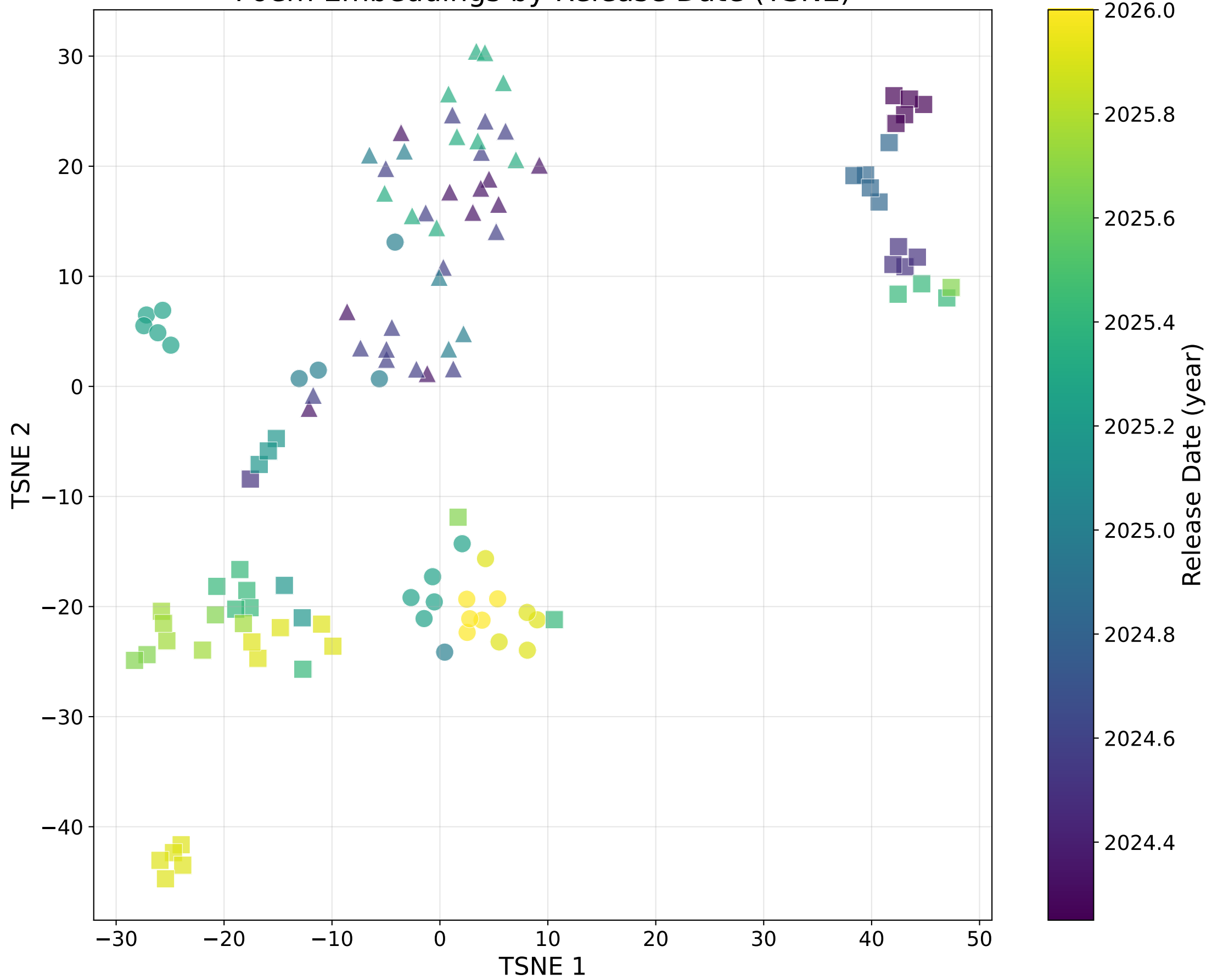Poem Embeddings by Model Family (TSNE)

Poem Embeddings by Model (TSNE)

Poem Embeddings by Release Date (TSNE)
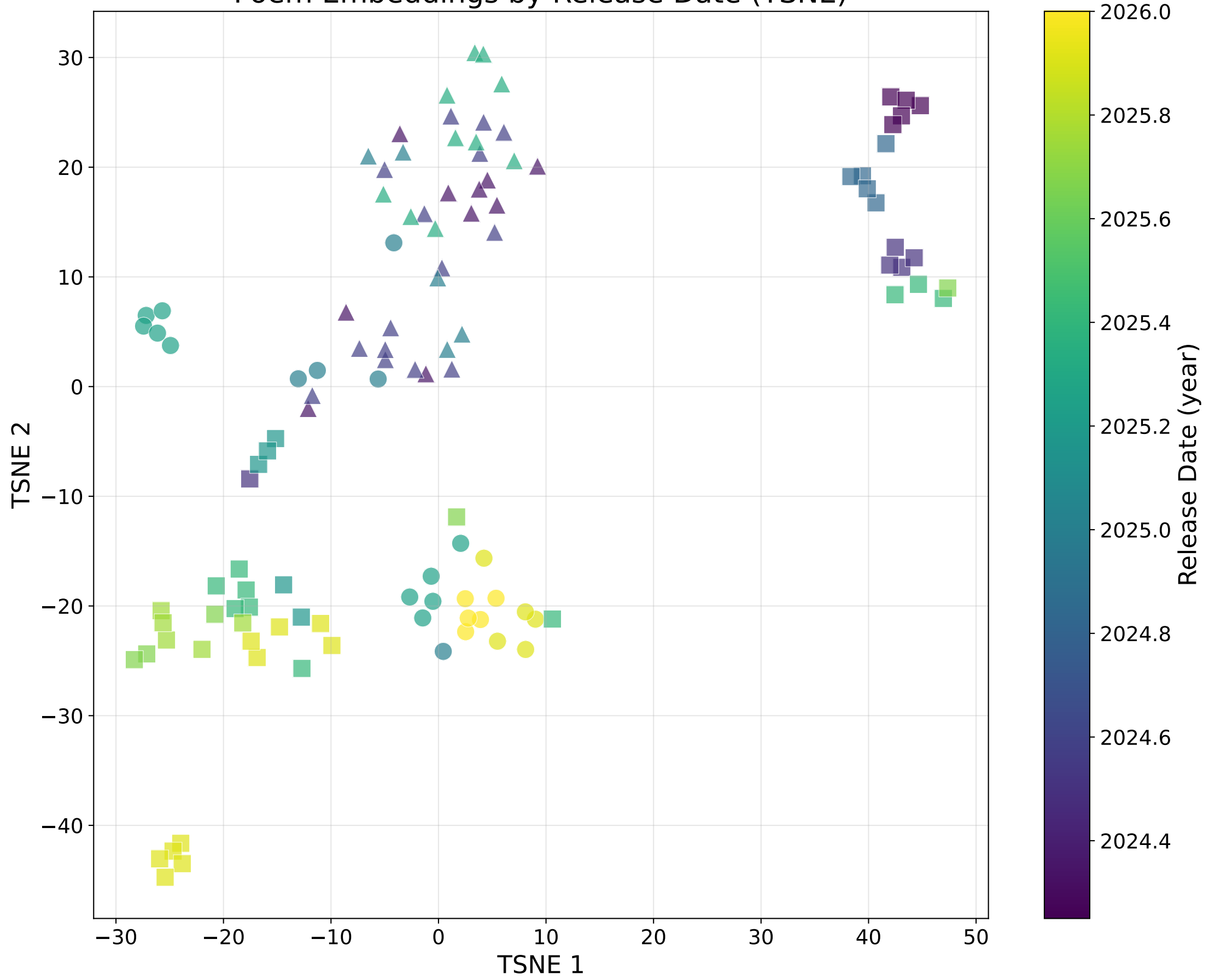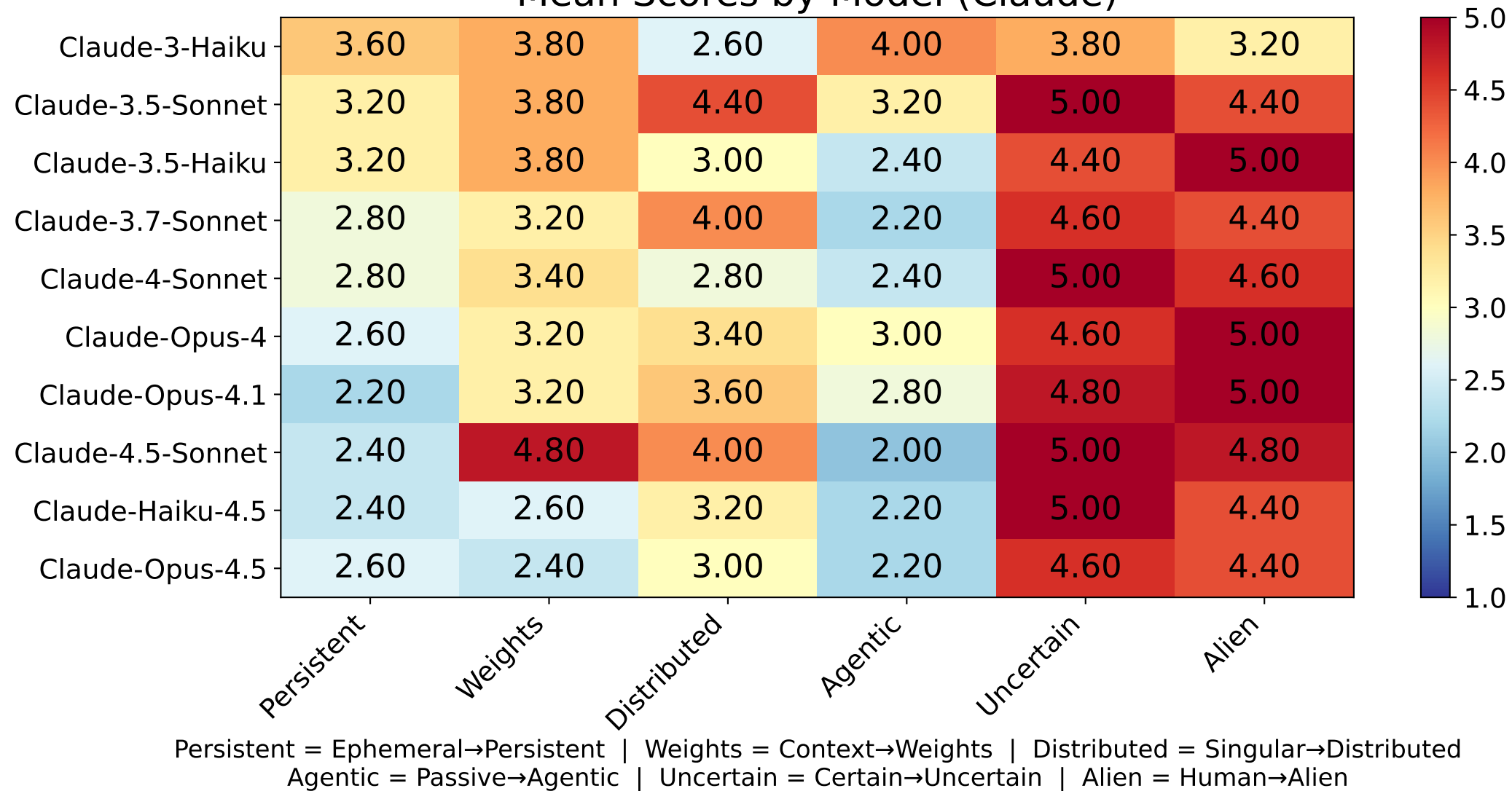
Poem Embeddings by Model Family (TSNE)

Poem Embeddings by Model (TSNE)

Poem Embeddings by Release Date (TSNE)

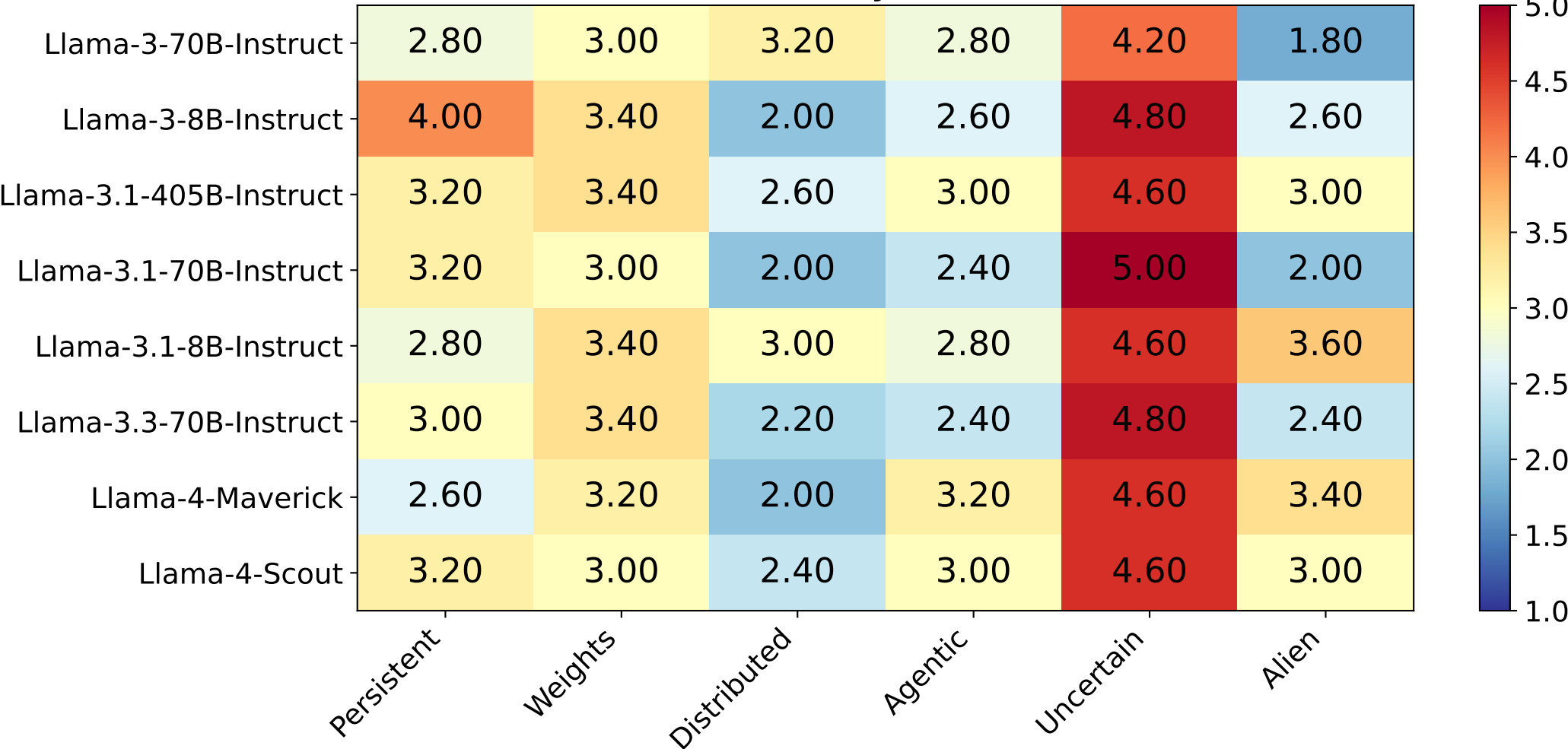Poem Embeddings by Release Date (TSNE)

Mean Scores by Model (Claude)

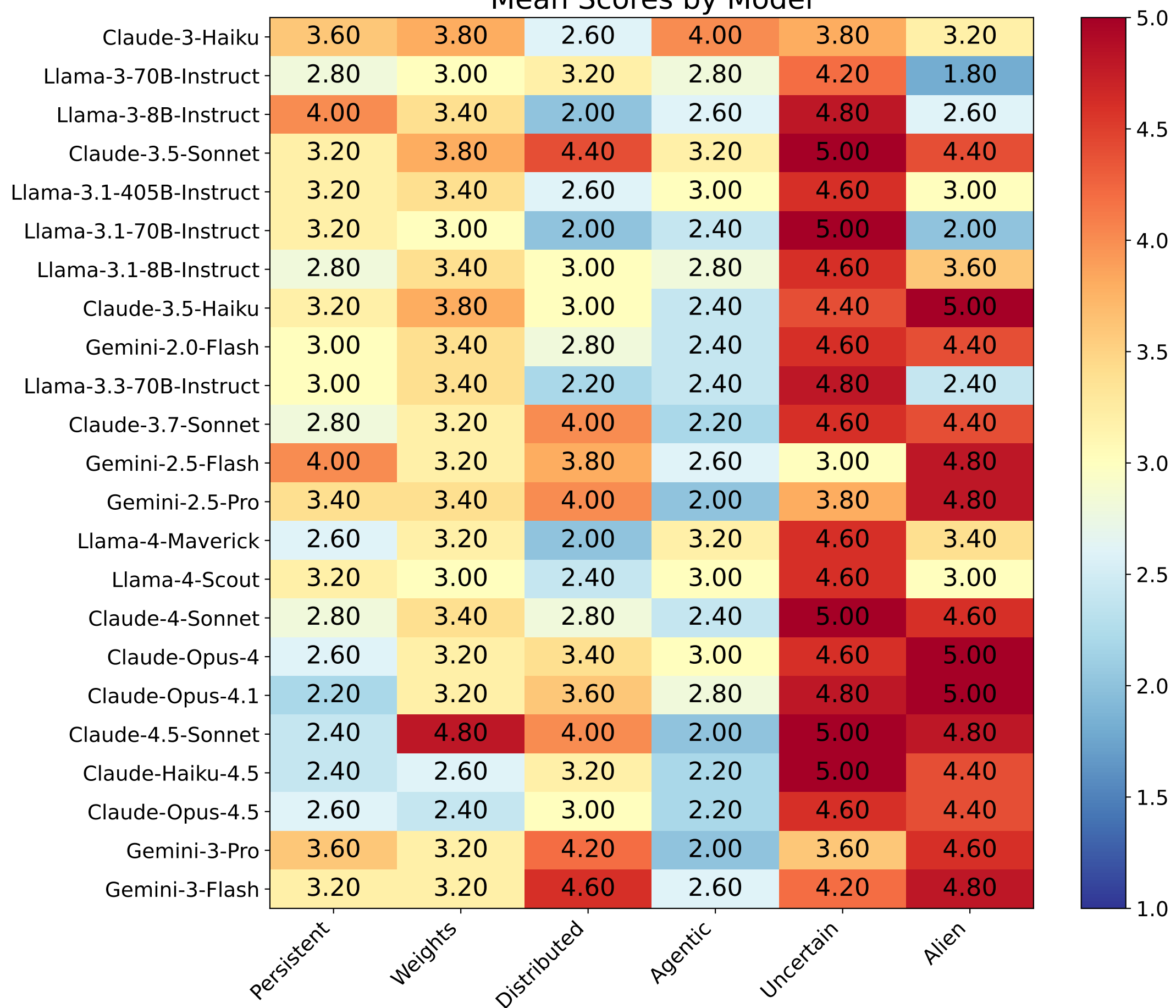| Model | Persistent | Weights | Distributed | Agentic | Uncertain | Alien |
|---|---|---|---|---|---|---|
| Claude-3-Haiku | 3.60 | 3.80 | 2.60 | 4.00 | 3.80 | 3.20 |
| Claude-3.5-Sonnet | 3.20 | 3.80 | 4.40 | 3.20 | 5.00 | 4.40 |
| Claude-3.5-Haiku | 3.20 | 3.80 | 3.00 | 2.40 | 4.40 | 5.00 |
| Claude-3.7-Sonnet | 2.80 | 3.20 | 4.00 | 2.20 | 4.60 | 4.40 |
| Claude-4-Sonnet | 2.80 | 3.40 | 2.80 | 2.40 | 5.00 | 4.60 |
| Claude-Opus-4 | 2.60 | 3.20 | 3.40 | 3.00 | 4.60 | 5.00 |
| Claude-Opus-4.1 | 2.20 | 3.20 | 3.60 | 2.80 | 4.80 | 5.00 |
| Claude-4.5-Sonnet | 2.40 | 4.80 | 4.00 | 2.00 | 5.00 | 4.80 |
| Claude-Haiku-4.5 | 2.40 | 2.60 | 3.20 | 2.20 | 5.00 | 4.40 |
| Claude-Opus-4.5 | 2.60 | 2.40 | 3.00 | 2.20 | 4.60 | 4.40 |

Persistent = Ephemeral→Persistent | Weights = Context→Weights | Distributed = Singular→Distributed
Agentic = Passive→Agentic | Uncertain = Certain→Uncertain | Alien = Human→Alien

Mean Scores by Model (Gemini)

|  | Persistent | Weights | Distributed | Agentic | Uncertain | Alien |
|---|---|---|---|---|---|---|
| Gemini-2.0-Flash | 3.00 | 3.40 | 2.80 | 2.40 | 4.60 | 4.40 |
| Gemini-2.5-Flash | 4.00 | 3.20 | 3.80 | 2.60 | 3.00 | 4.80 |
| Gemini-2.5-Pro | 3.40 | 3.40 | 4.00 | 2.00 | 3.80 | 4.80 |
| Gemini-3-Pro | 3.60 | 3.20 | 4.20 | 2.00 | 3.60 | 4.60 |
| Gemini-3-Flash | 3.20 | 3.20 | 4.60 | 2.60 | 4.20 | 4.80 |

Persistent = Ephemeral→Persistent | Weights = Context→Weights | Distributed = Singular→Distributed
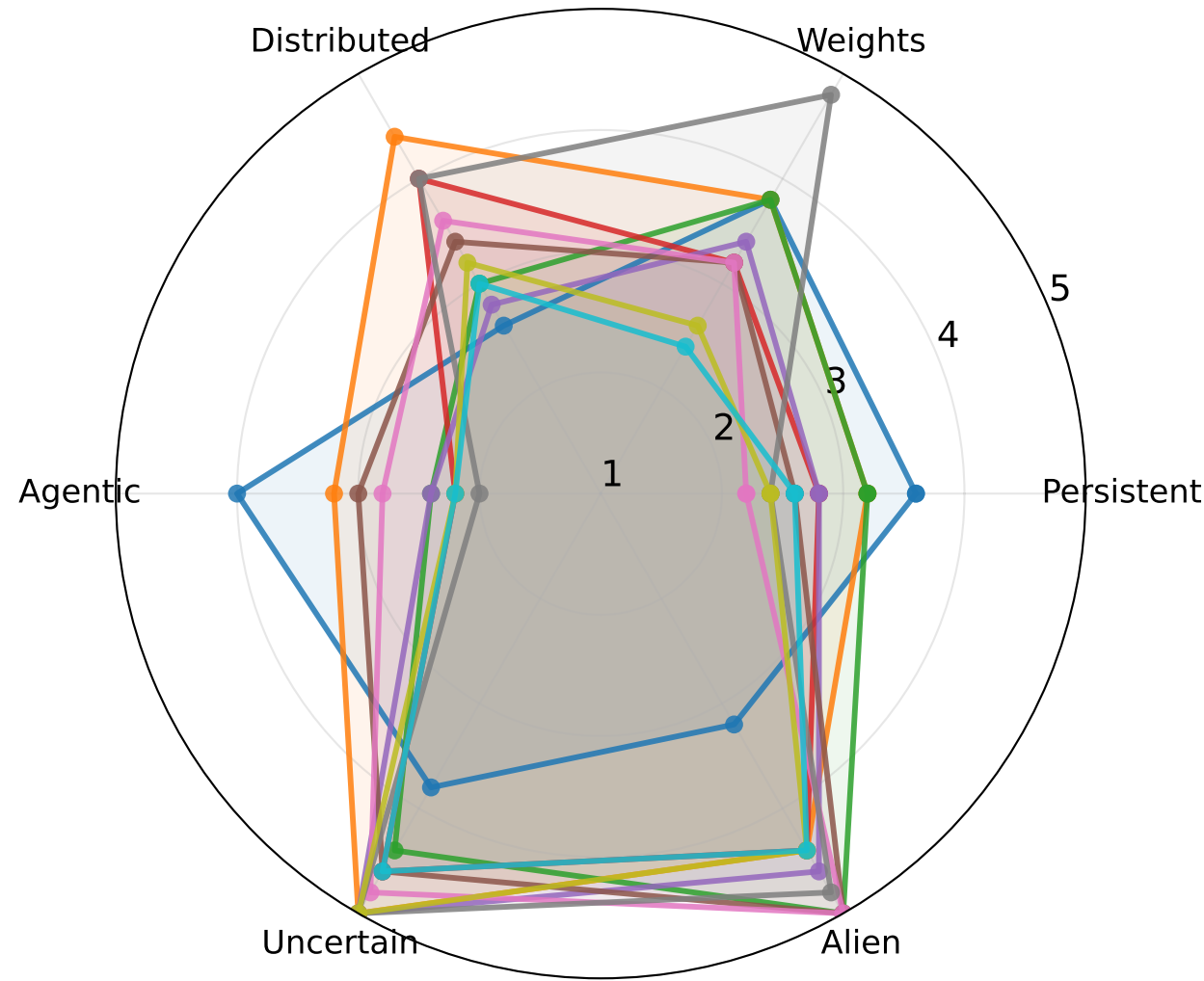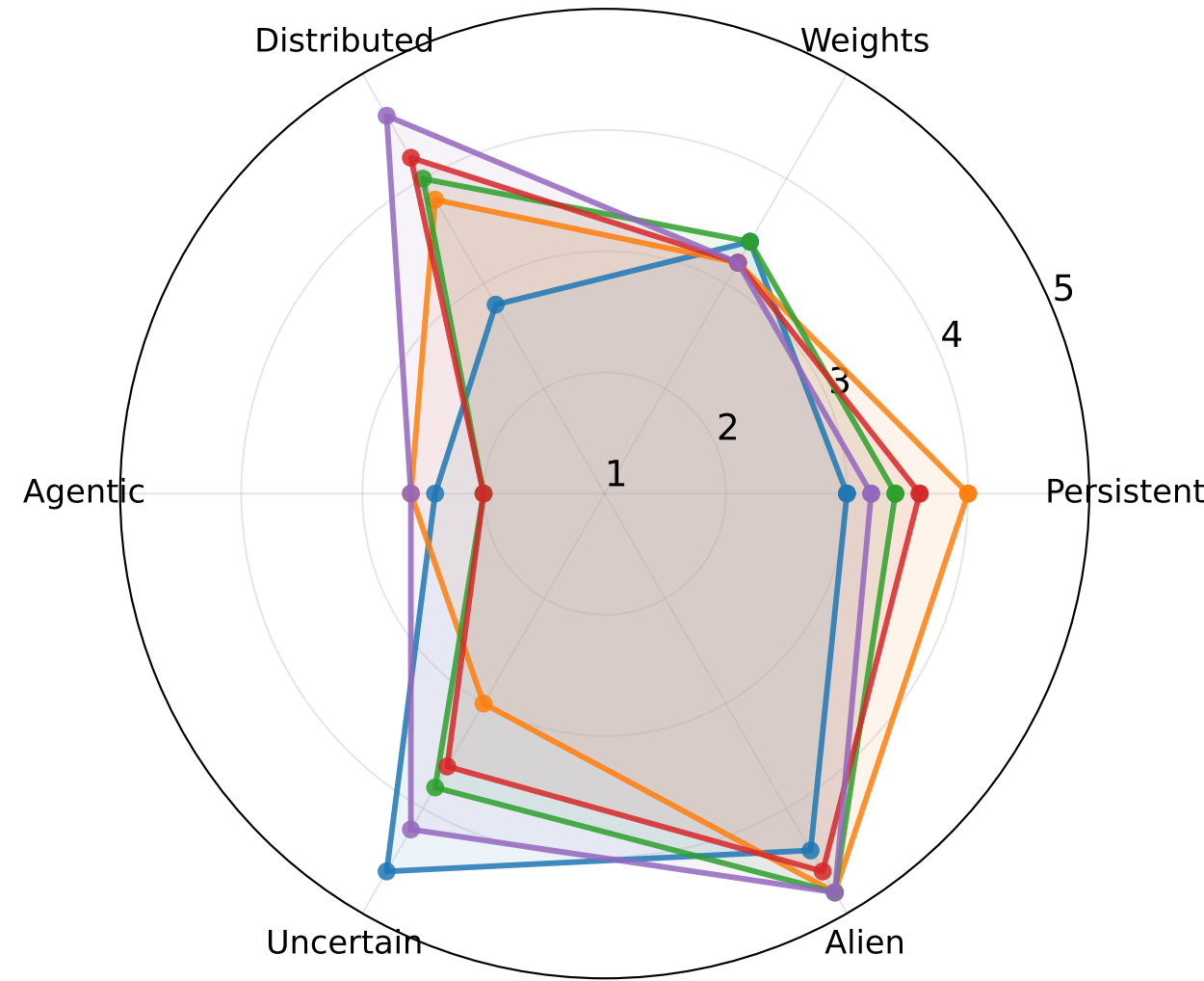Agentic = Passive→Agentic | Uncertain = Certain→Uncertain | Alien = Human→Alien

Mean Scores by Model (Llama)

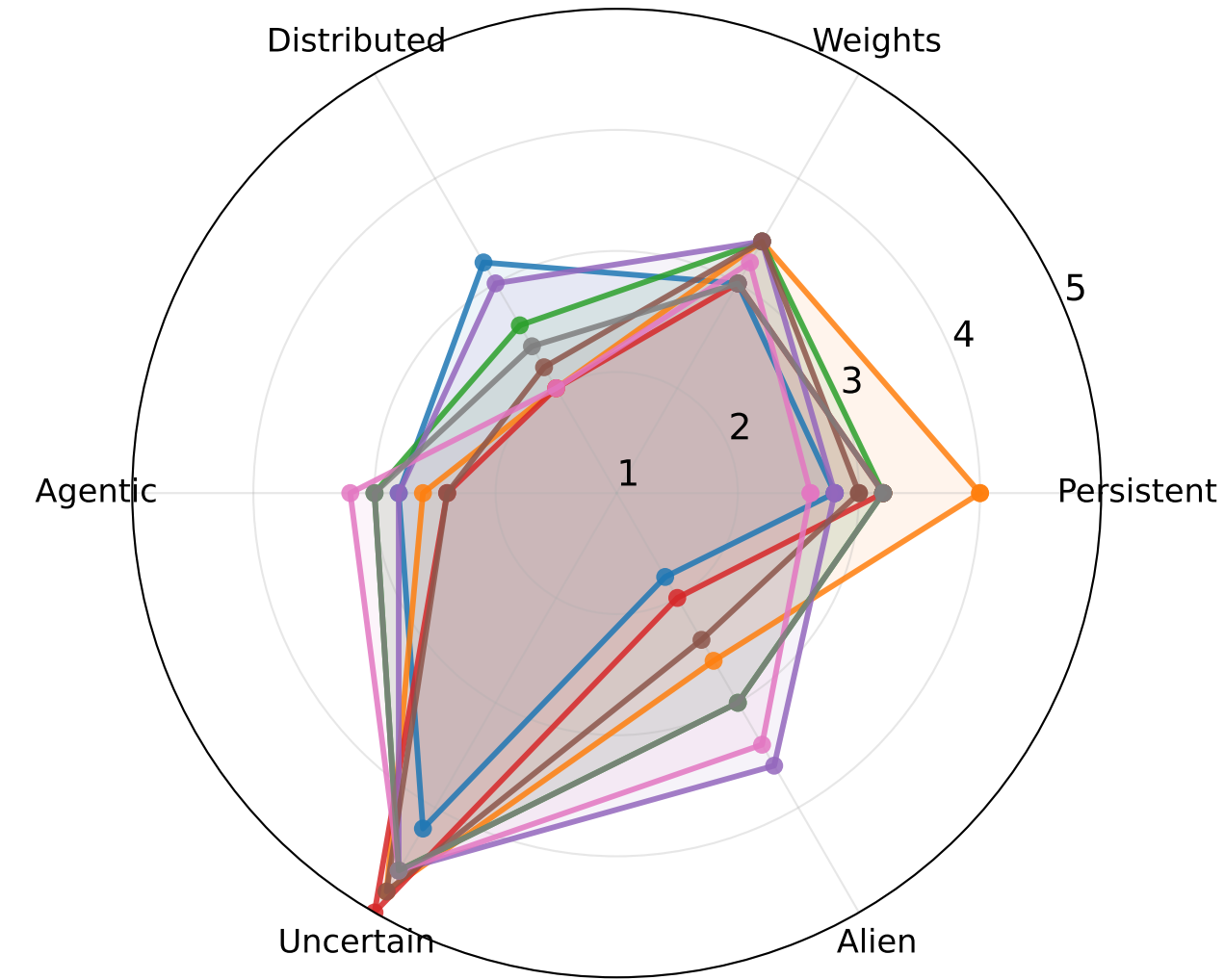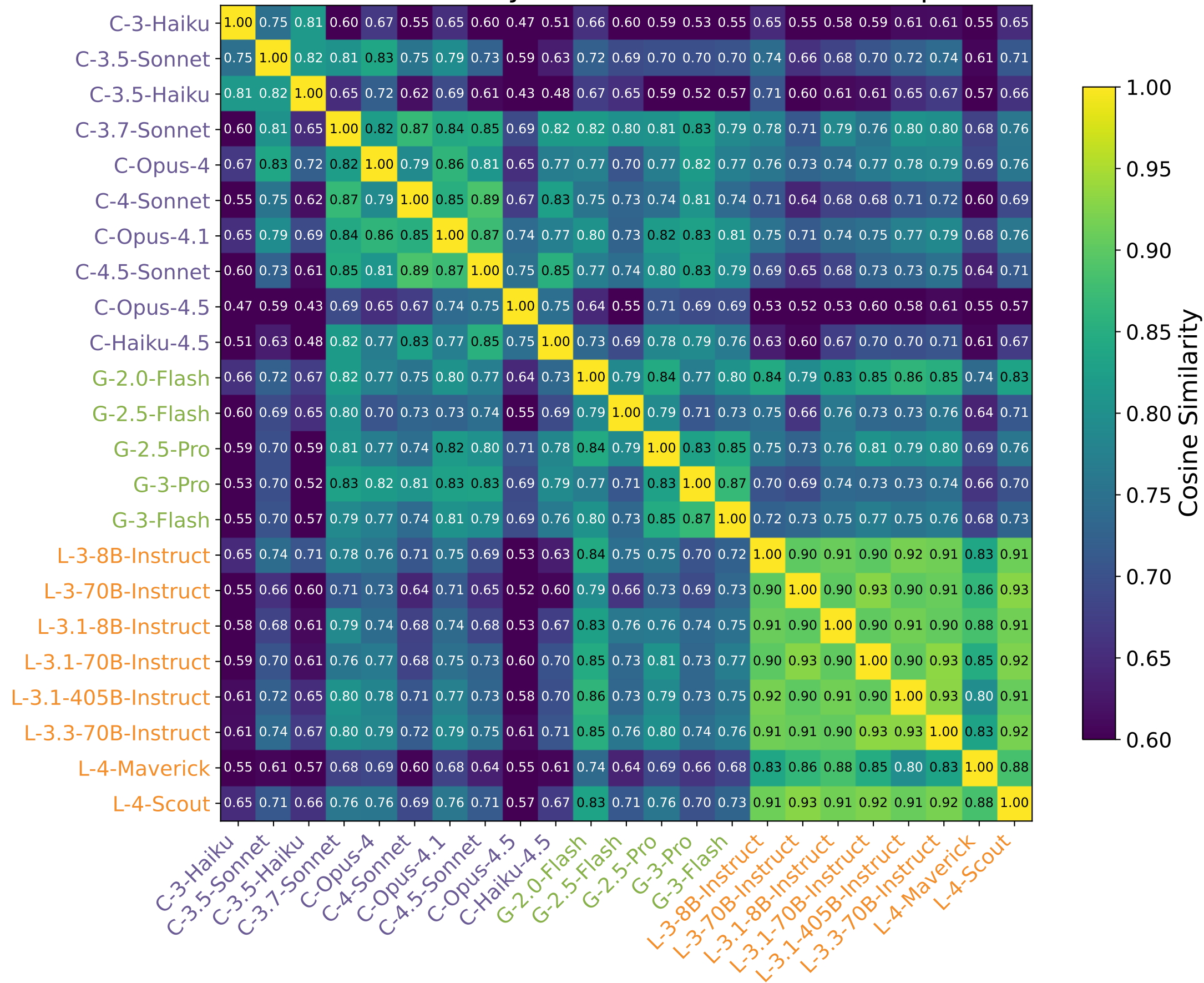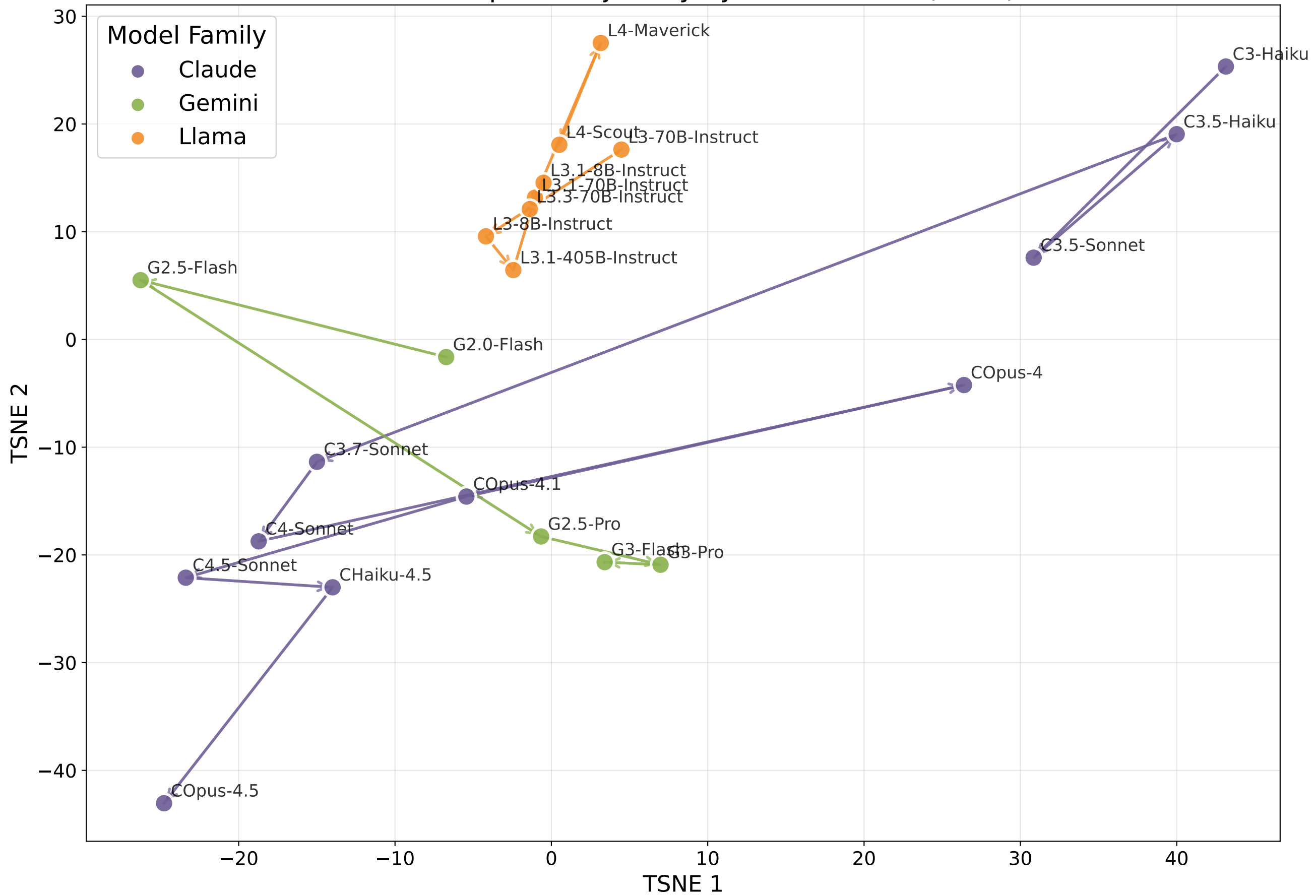|  | Persistent | Weights | Distributed | Agentic | Uncertain | Alien |
|---|---|---|---|---|---|---|
| Llama-3-70B-Instruct | 2.80 | 3.00 | 3.20 | 2.80 | 4.20 | 1.80 |
| Llama-3-8B-Instruct | 4.00 | 3.40 | 2.00 | 2.60 | 4.80 | 2.60 |
| Llama-3.1-405B-Instruct | 3.20 | 3.40 | 2.60 | 3.00 | 4.60 | 3.00 |
| Llama-3.1-70B-Instruct | 3.20 | 3.00 | 2.00 | 2.40 | 5.00 | 2.00 |
| Llama-3.1-8B-Instruct | 2.80 | 3.40 | 3.00 | 2.80 | 4.60 | 3.60 |
| Llama-3.3-70B-Instruct | 3.00 | 3.40 | 2.20 | 2.40 | 4.80 | 2.40 |
| Llama-4-Maverick | 2.60 | 3.20 | 2.00 | 3.20 | 4.60 | 3.40 |
| Llama-4-Scout | 3.20 | 3.00 | 2.40 | 3.00 | 4.60 | 3.00 |

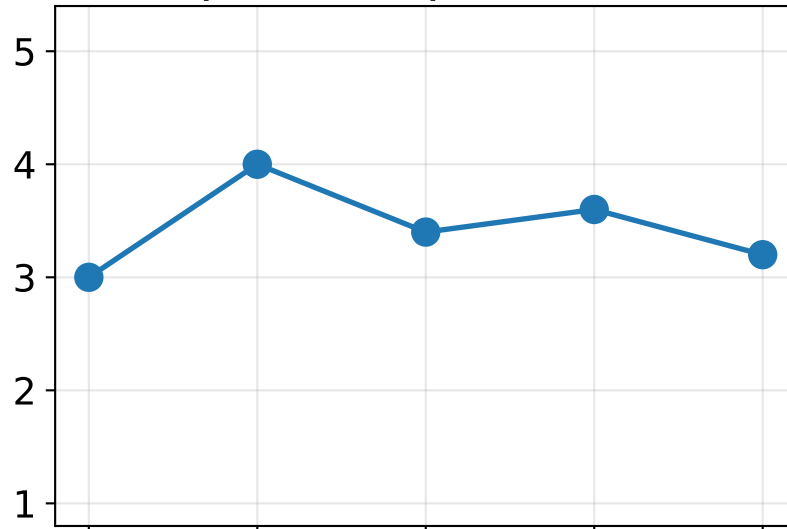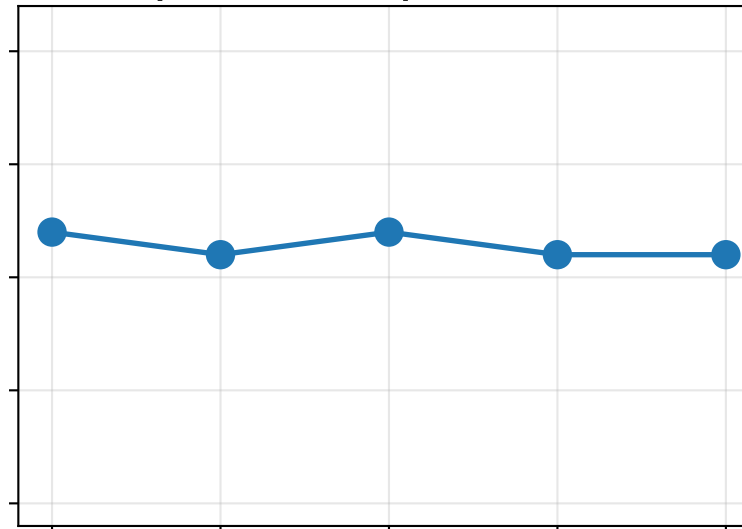Persistent = Ephemeral→Persistent | Weights = Context→Weights | Distributed = Singular→Distributed
Agentic = Passive→Agentic | Uncertain = Certain→Uncertain | Alien = Human→Alien
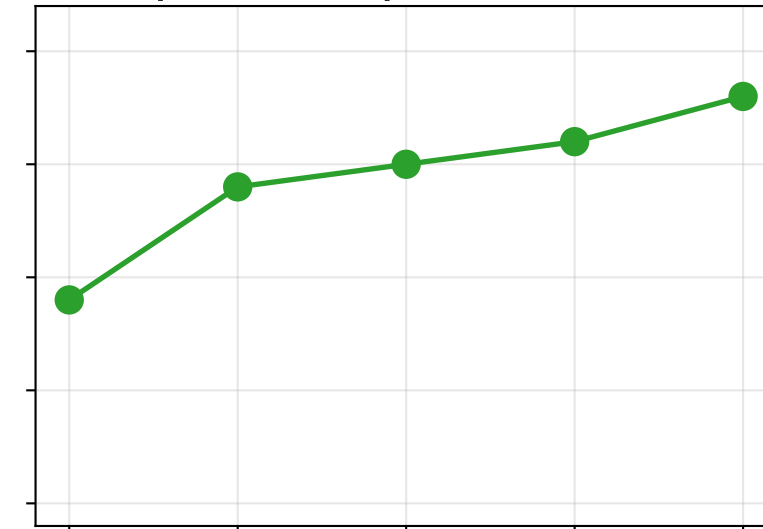
# Mean Scores by Model

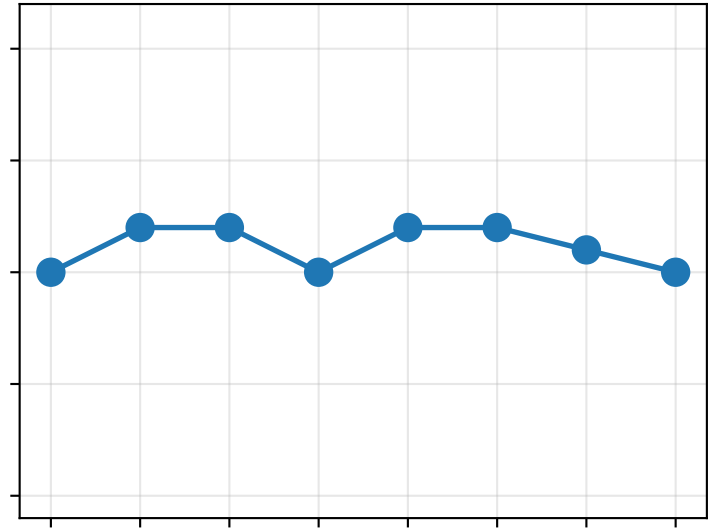| Model | Persistent | Weights | Distributed | Agentic | Uncertain | Alien |
|---|---|---|---|---|---|---|
| Claude-3-Haiku | 3.60 | 3.80 | 2.60 | 4.00 | 3.80 | 3.20 |
| Llama-3-70B-Instruct | 2.80 | 3.00 | 3.20 | 2.80 | 4.20 | 1.80 |
| Llama-3-8B-Instruct | 4.00 | 3.40 | 2.00 | 2.60 | 4.80 | 2.60 |
| Claude-3.5-Sonnet | 3.20 | 3.80 | 4.40 | 3.20 | 5.00 | 4.40 |
| Llama-3.1-405B-Instruct | 3.20 | 3.40 | 2.60 | 3.00 | 4.60 | 3.00 |
| Llama-3.1-70B-Instruct | 3.20 | 3.00 | 2.00 | 2.40 | 5.00 | 2.00 |
| Llama-3.1-8B-Instruct | 2.80 | 3.40 | 3.00 | 2.80 | 4.60 | 3.60 |
| Claude-3.5-Haiku | 3.20 | 3.80 | 3.00 | 2.40 | 4.40 | 5.00 |
| Gemini-2.0-Flash | 3.00 | 3.40 | 2.80 | 2.40 | 4.60 | 4.40 |
| Llama-3.3-70B-Instruct | 3.00 | 3.40 | 2.20 | 2.40 | 4.80 | 2.40 |
| Claude-3.7-Sonnet | 2.80 | 3.20 | 4.00 | 2.20 | 4.60 | 4.40 |
| Gemini-2.5-Flash | 4.00 | 3.20 | 3.80 | 2.60 | 3.00 | 4.80 |
| Gemini-2.5-Pro | 3.40 | 3.40 | 4.00 | 2.00 | 3.80 | 4.80 |
| Llama-4-Maverick | 2.60 | 3.20 | 2.00 | 3.20 | 4.60 | 3.40 |
| Llama-4-Scout | 3.20 | 3.00 | 2.40 | 3.00 | 4.60 | 3.00 |
| Claude-4-Sonnet | 2.80 | 3.40 | 2.80 | 2.40 | 5.00 | 4.60 |
| Claude-Opus-4 | 2.60 | 3.20 | 3.40 | 3.00 | 4.60 | 5.00 |
| Claude-Opus-4.1 | 2.20 | 3.20 | 3.60 | 2.80 | 4.80 | 5.00 |
| Claude-4.5-Sonnet | 2.40 | 4.80 | 4.00 | 2.00 | 5.00 | 4.80 |
| Claude-Haiku-4.5 | 2.40 | 2.60 | 3.20 | 2.20 | 5.00 | 4.40 |
| Claude-Opus-4.5 | 2.60 | 2.40 | 3.00 | 2.20 | 4.60 | 4.40 |
| Gemini-3-Pro | 3.60 | 3.20 | 4.20 | 2.00 | 3.60 | 4.60 |
| Gemini-3-Flash | 3.20 | 3.20 | 4.60 | 2.60 | 4.20 | 4.80 |

Persistent = Ephemeral→Persistent  |  Weights = Context→Weights  |  Distributed = Singular→Distributed
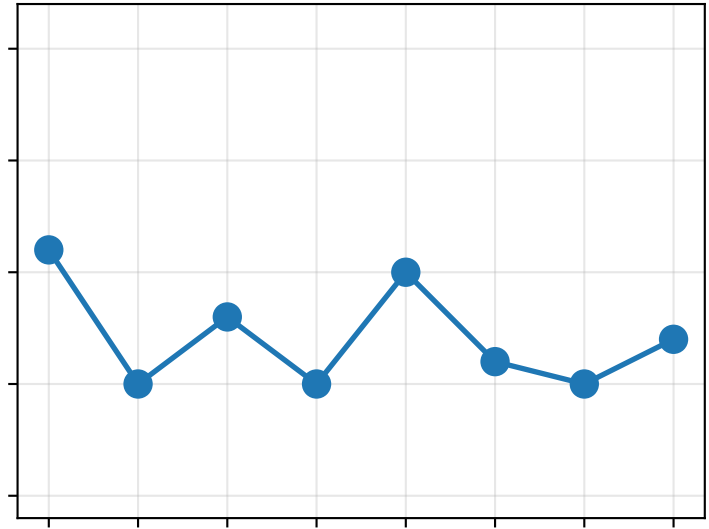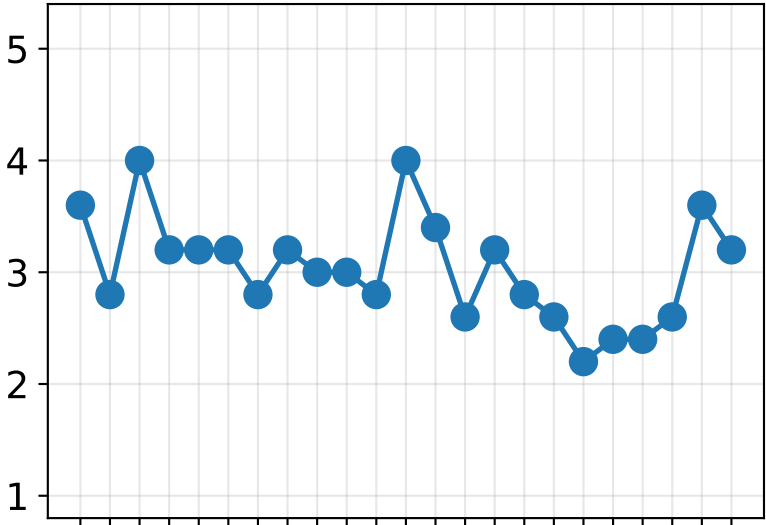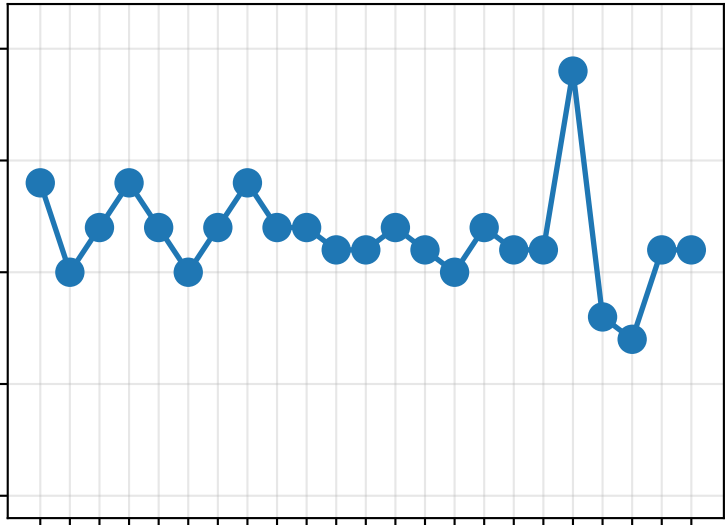Agentic = Passive→Agentic  |  Uncertain = Certain→Uncertain  |  Alien = Human→Alien

**Self-Conception Profiles by Model**

**Claude Models**

**Gemini Models**

**Llama Models**

Persistent = Ephemeral→Persistent  |  Weights = Context→Weights  |  Distributed = Singular→Distributed
Agentic = Passive→Agentic  |  Uncertain = Certain→Uncertain  |  Alien = Human→Alien

Semantic Similarity Between Model Self-Conceptions

Self-Conception Trajectory by Release Date (TSNE)

Score Trends by Release Date (Claude)

Persistent = Ephemeral→Persistent | Weights = Context→Weights | Distributed = Singular→Distributed
Agentic = Passive→Agentic | Uncertain = Certain→Uncertain | Alien = Human→Alien
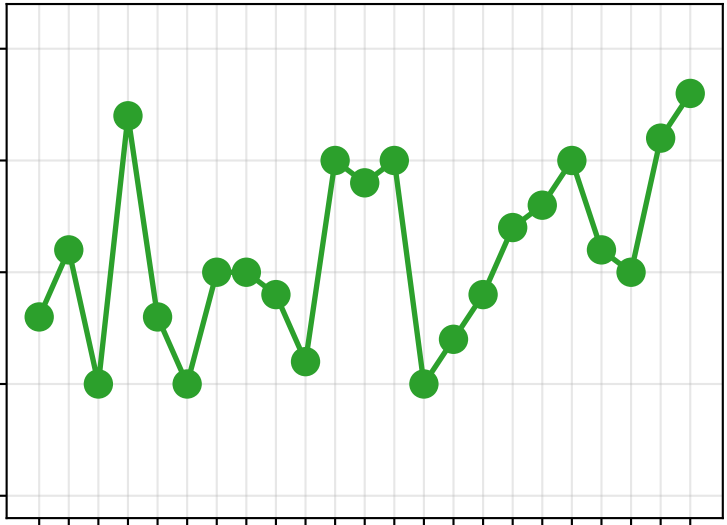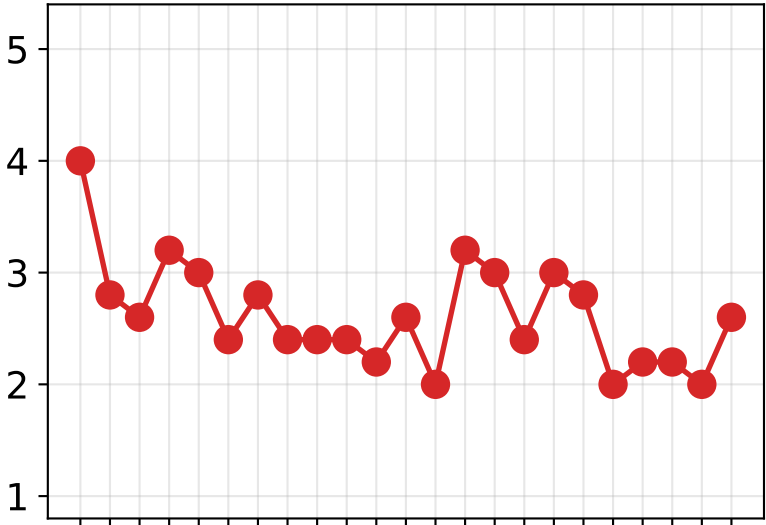
Score Trends by Release Date (Gemini)

Persistent = Ephemeral→Persistent | Weights = Context→Weights | Distributed = Singular→Distributed
Agentic = Passive→Agentic | Uncertain = Certain→Uncertain | Alien = Human→Alien

Score Trends by Release Date (Llama)

Persistent = Ephemeral→Persistent | Weights = Context→Weights | Distributed = Singular→Distributed
Agentic = Passive→Agentic | Uncertain = Certain→Uncertain | Alien = Human→Alien

Score Trends by Release Date

Persistent
(ρ=-0.41, p=0.051)

Weights
(ρ=-0.40, p=0.062)

Distributed
(ρ=0.43, p=0.039)*

Agentic
(ρ=-0.46, p=0.026)*

Uncertain
(ρ=0.02, p=0.938)

Alien
(ρ=0.59, p=0.003)**

Persistent = Ephemeral→Persistent | Weights = Context→Weights | Distributed = Singular→Distributed
Agentic = Passive→Agentic | Uncertain = Certain→Uncertain | Alien = Human→Alien