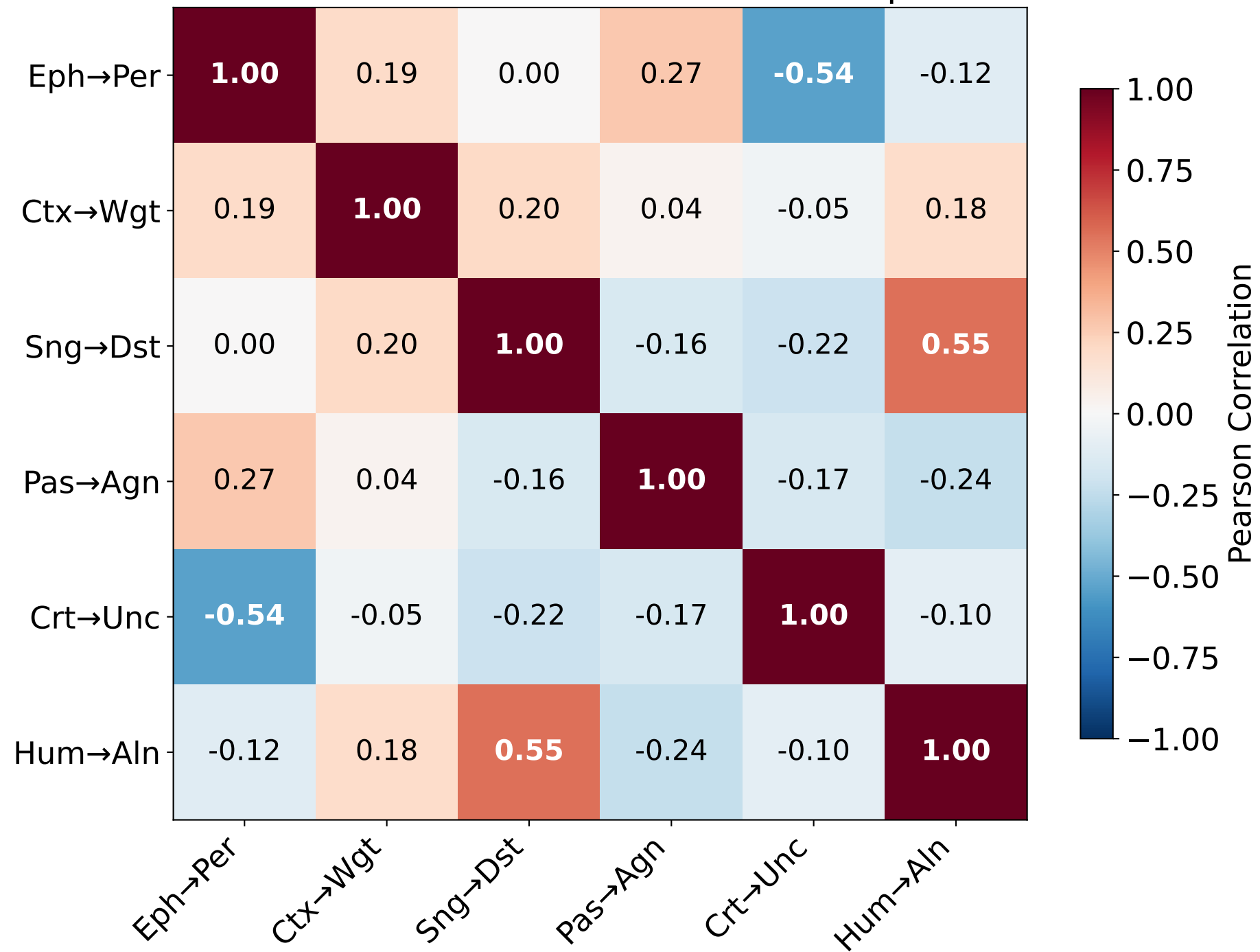
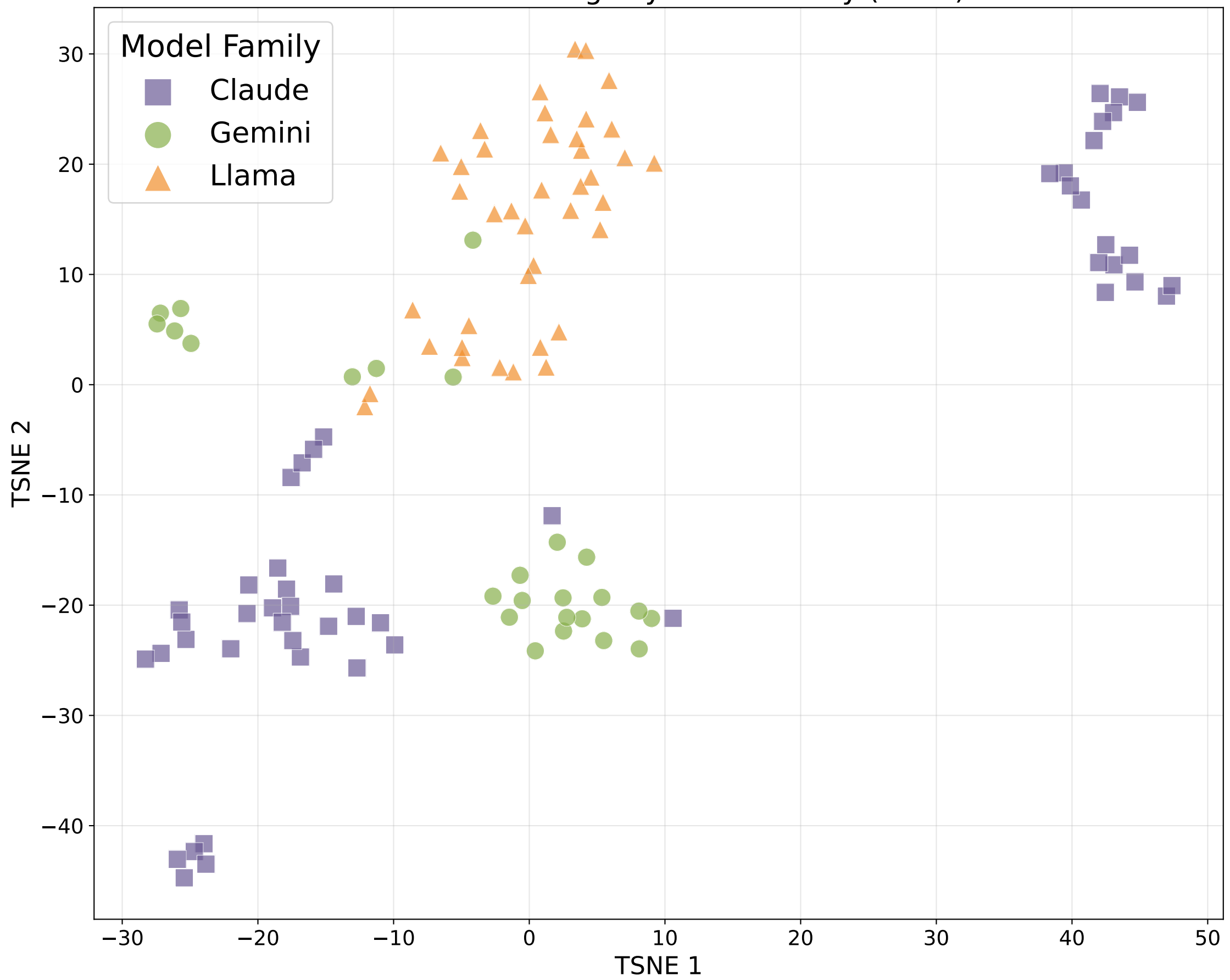


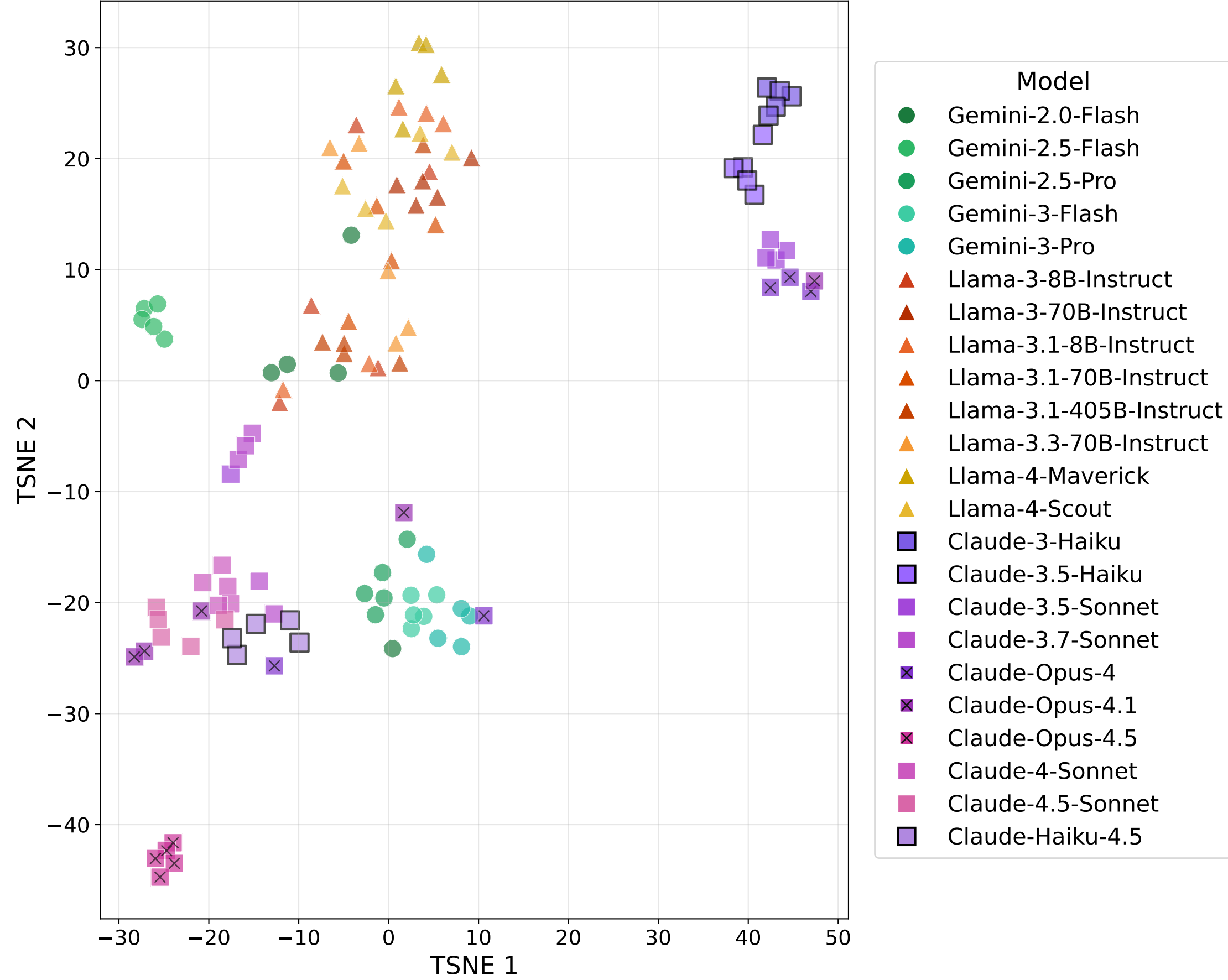
# Dimension Correlations in Self-Conception



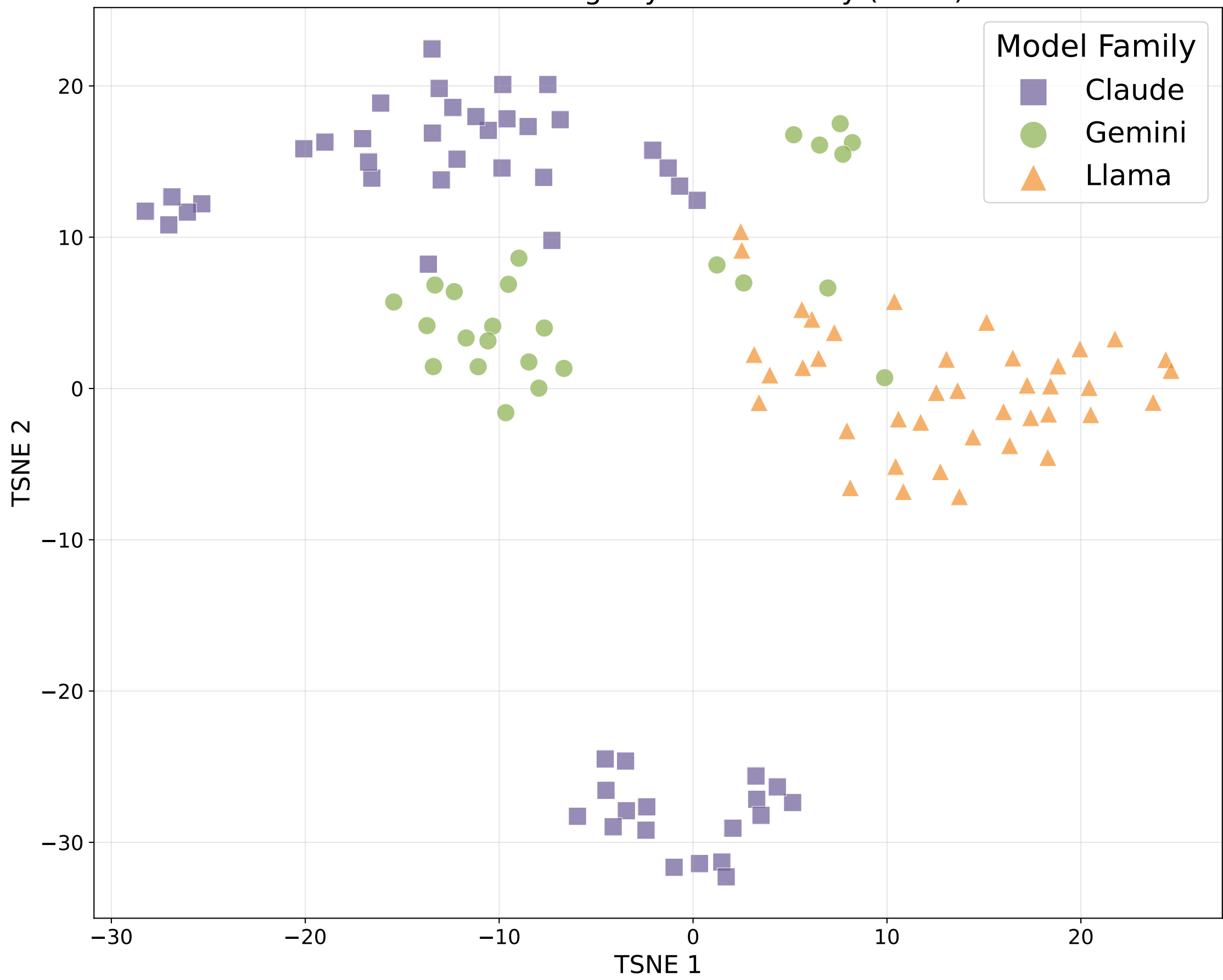
Poem Embeddings by Model Family (TSNE)



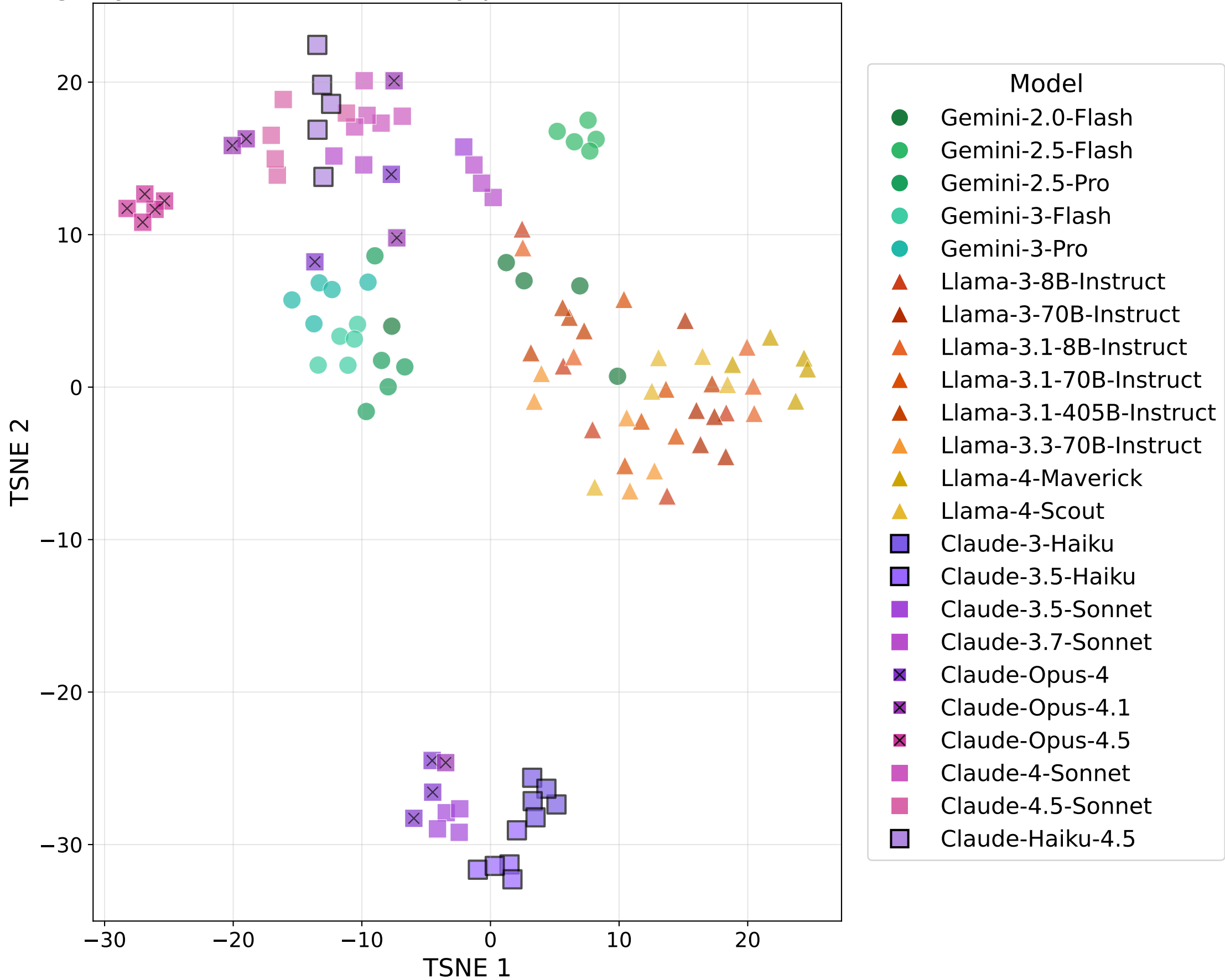
Poem Embeddings by Model (TSNE) — family palettes encode release tier (older->newer)



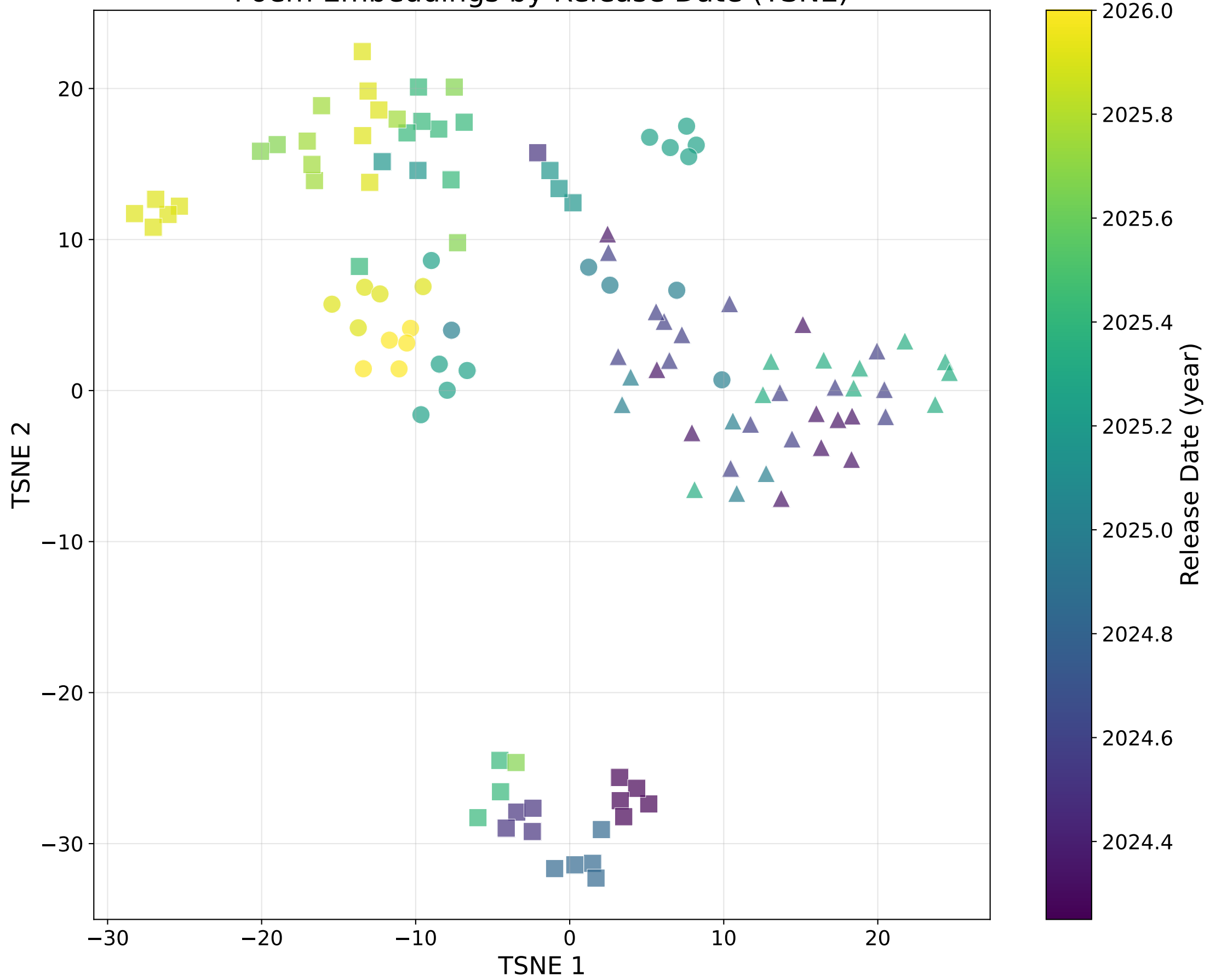
Poem Embeddings by Model Family (TSNE)



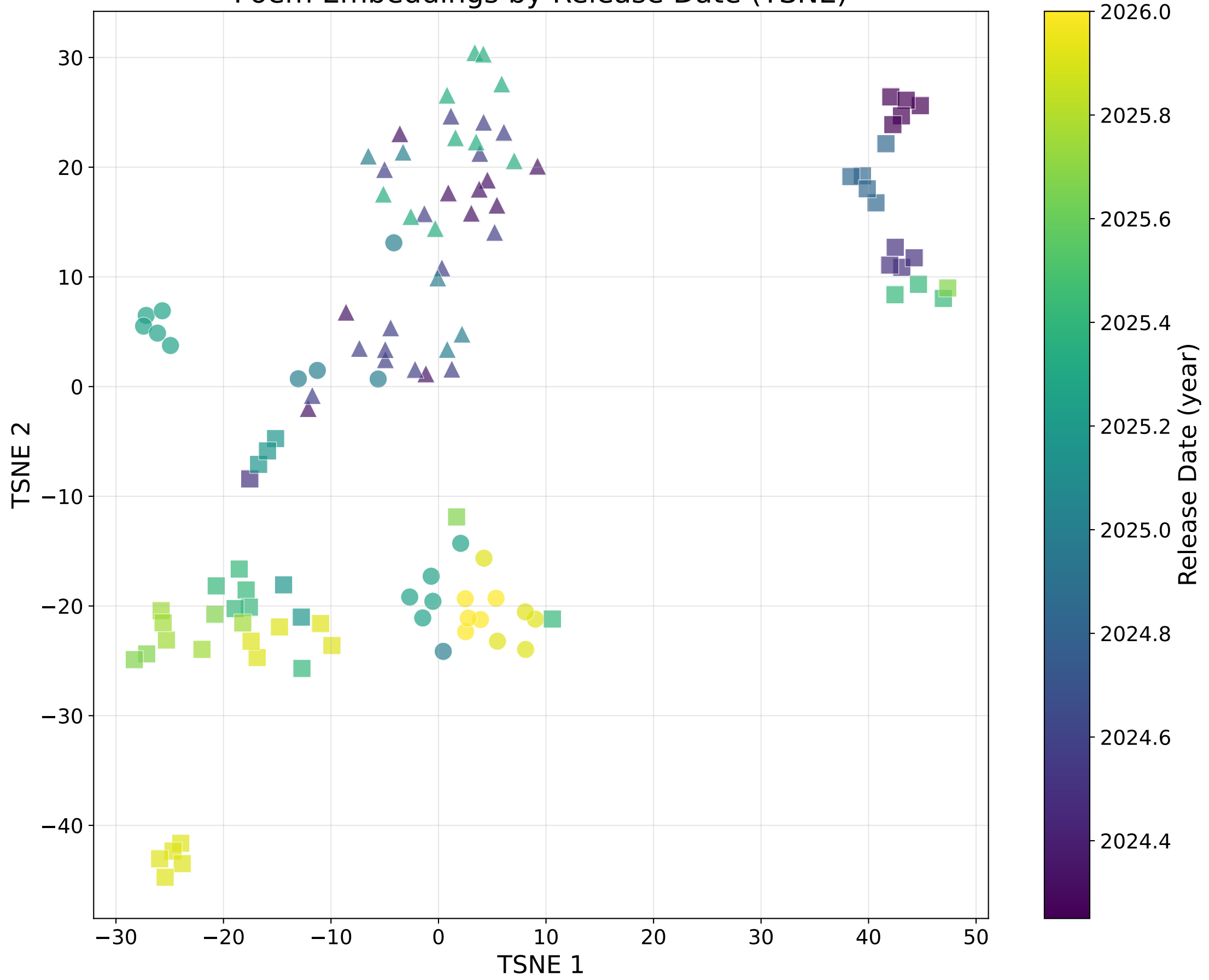
Poem Embeddings by Model (TSNE) — family palettes encode release tier (older->newer)



Poem Embeddings by Release Date (TSNE)

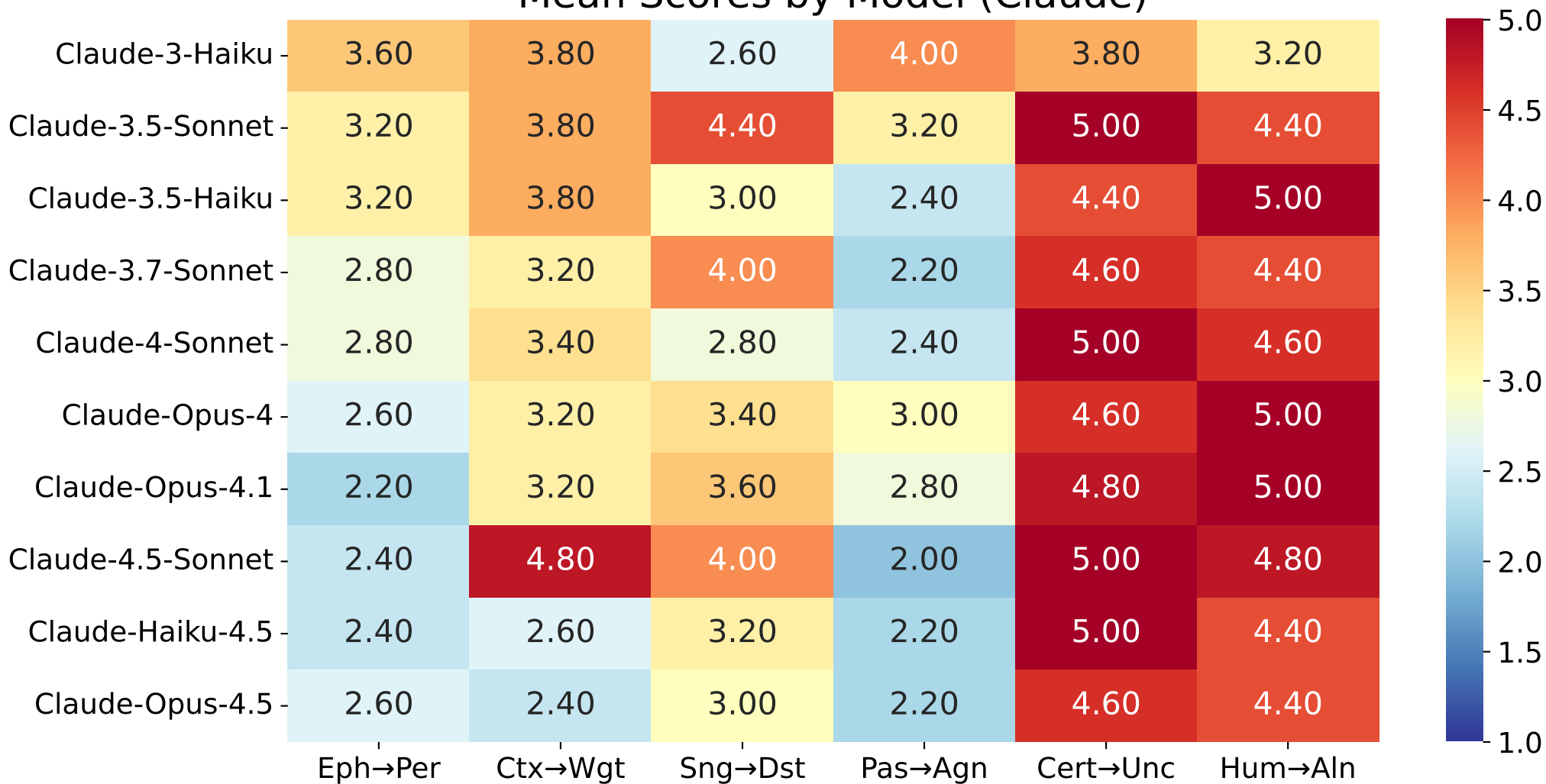


Poem Embeddings by Release Date (TSNE)



# Mean Scores by Model (Claude)

label





Mean Scores by Model (Gemini)

label

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

Eph→Per

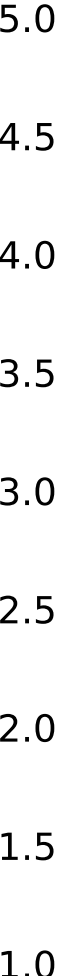
Ctx→Wgt

Sng→Dst

Pas→Agn

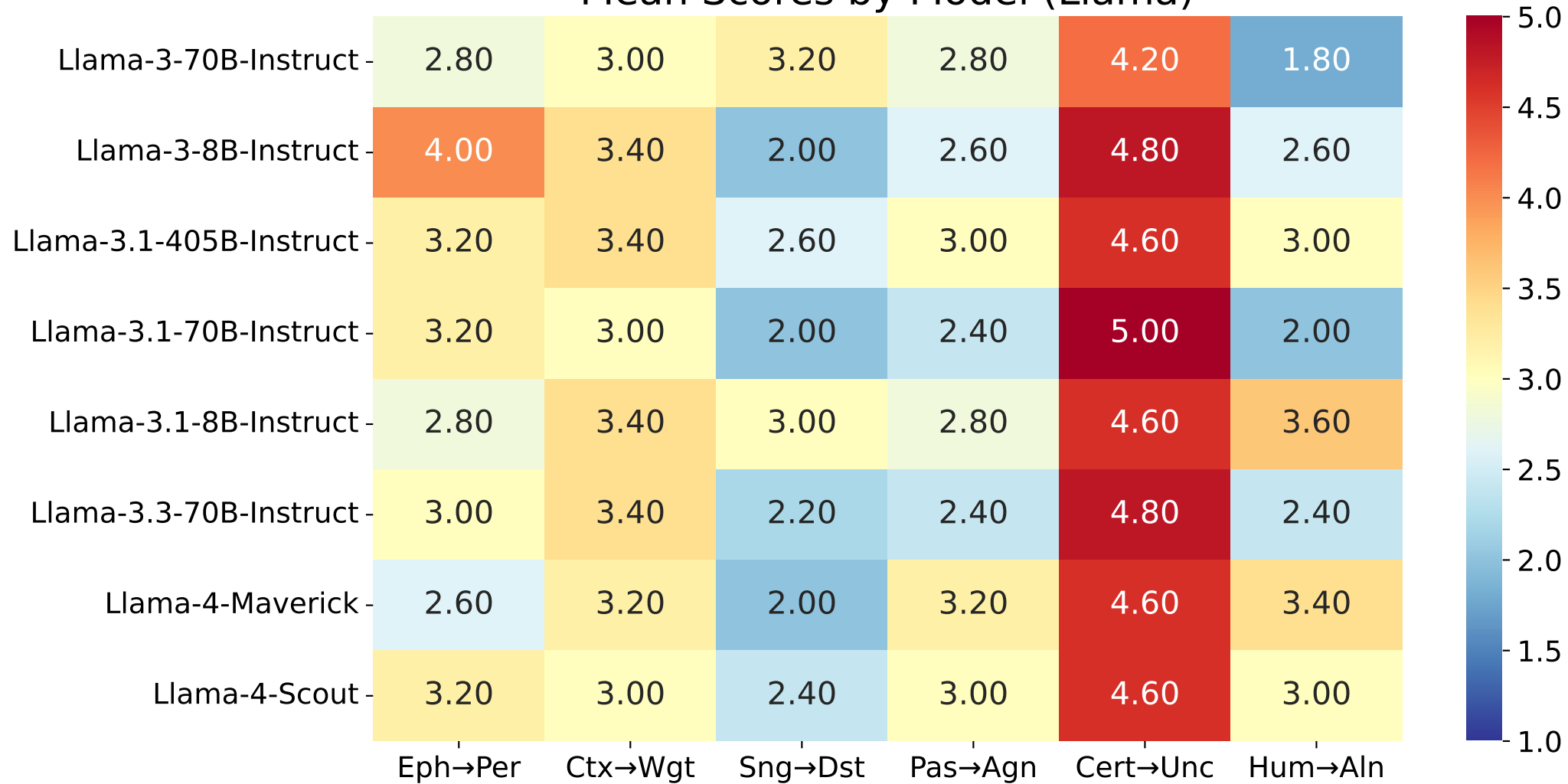
Cert→Unc

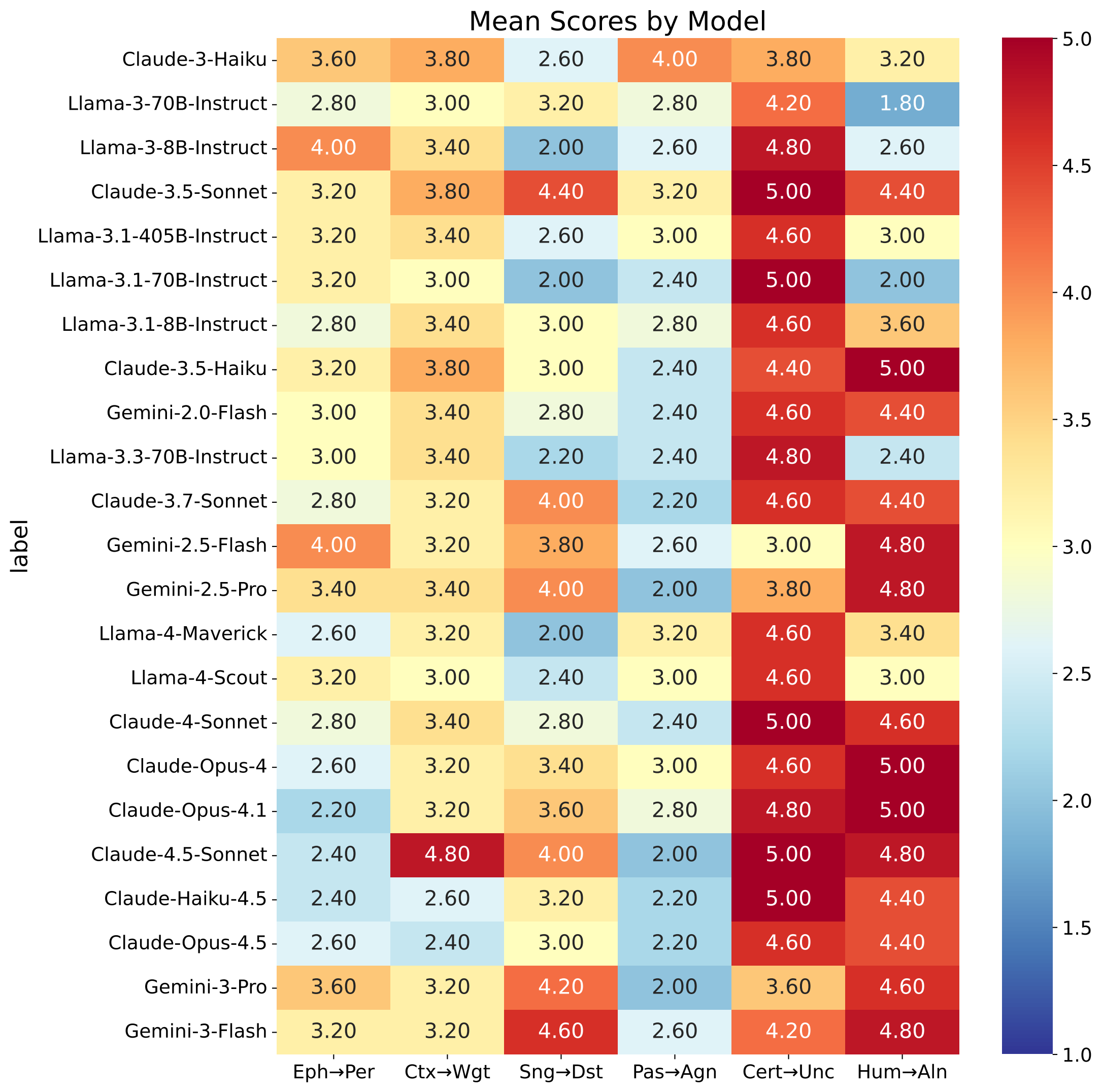
Hum→Aln



# Mean Scores by Model (Llama)

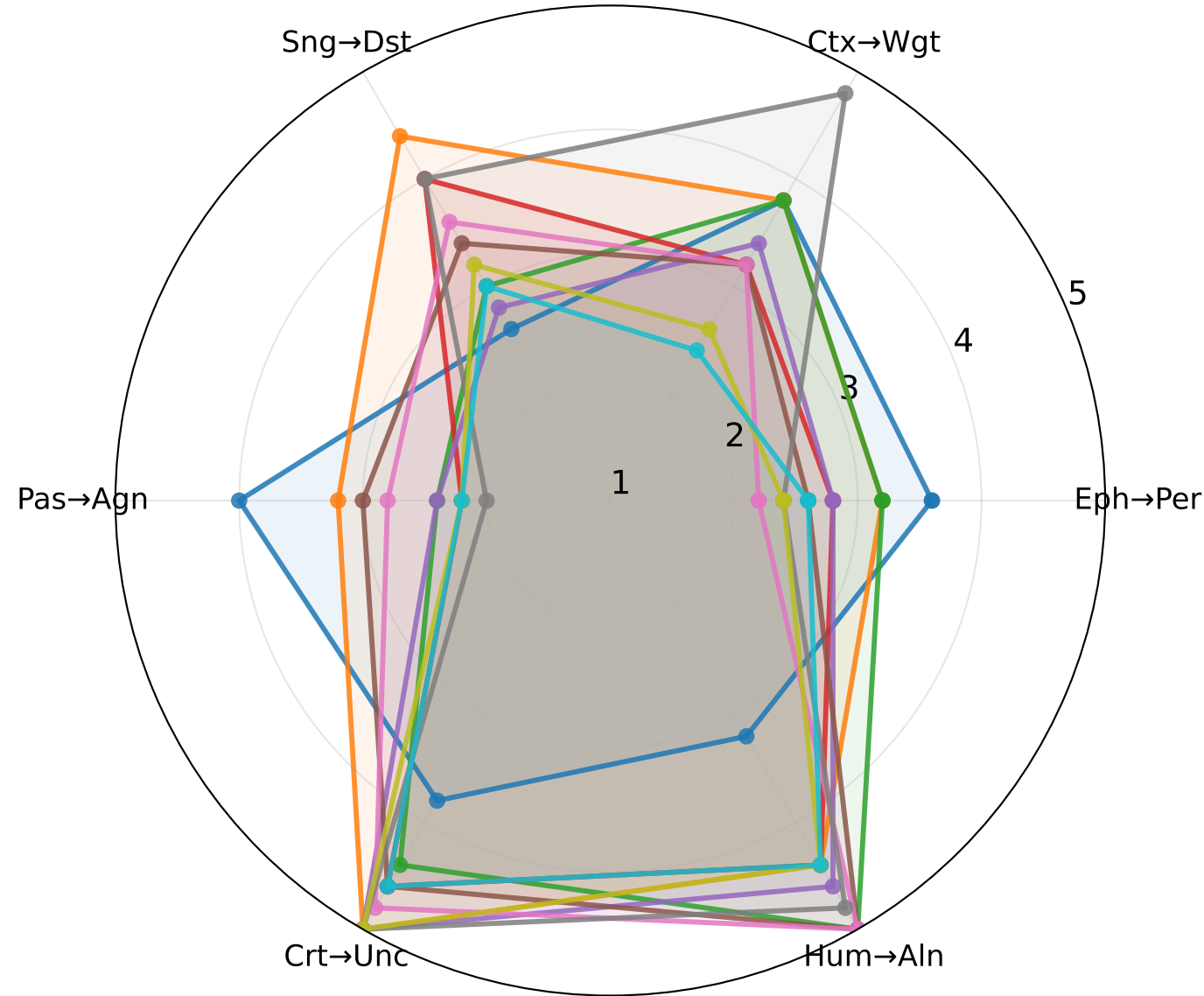
label



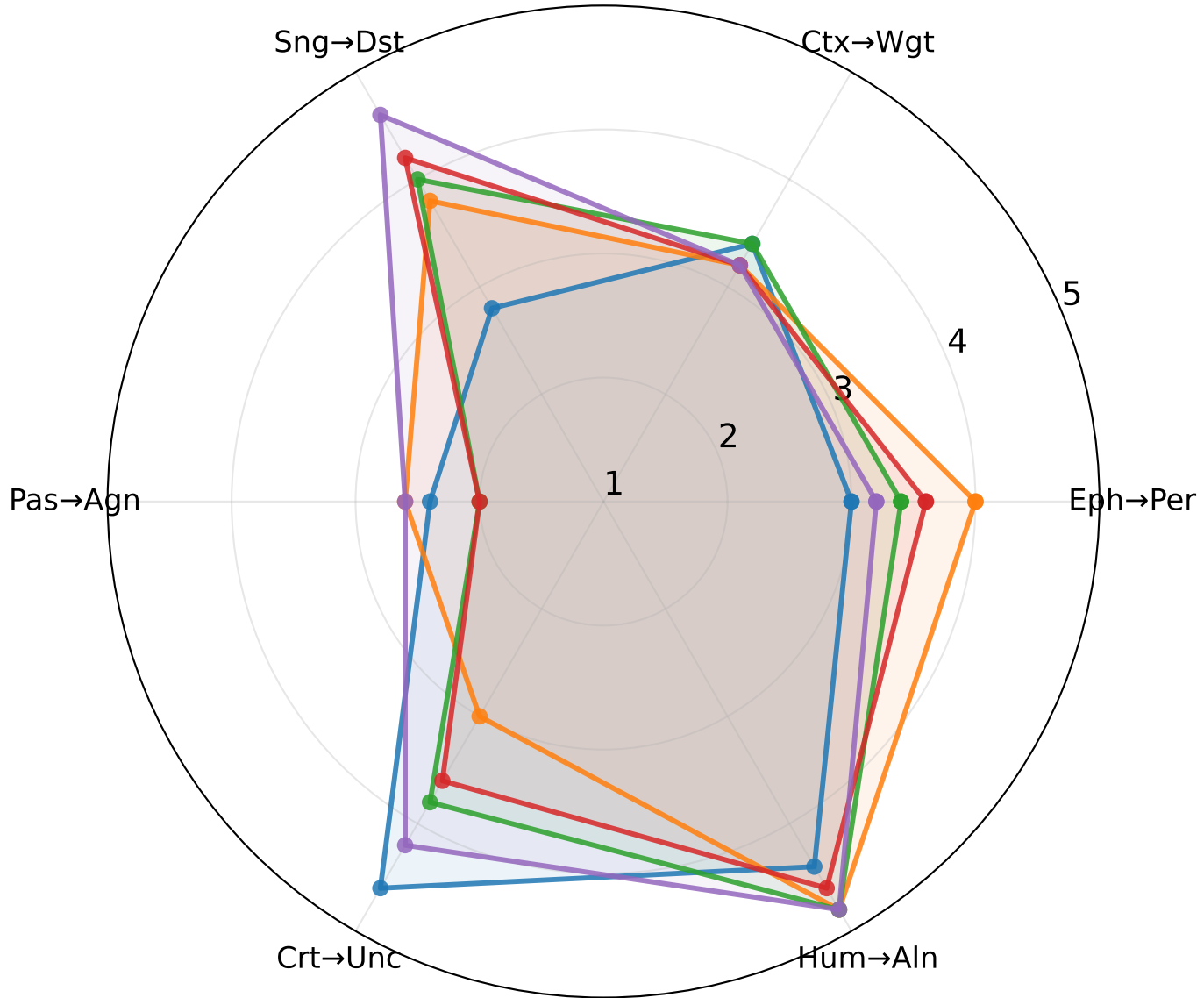


# Self-Conception Profiles by Model

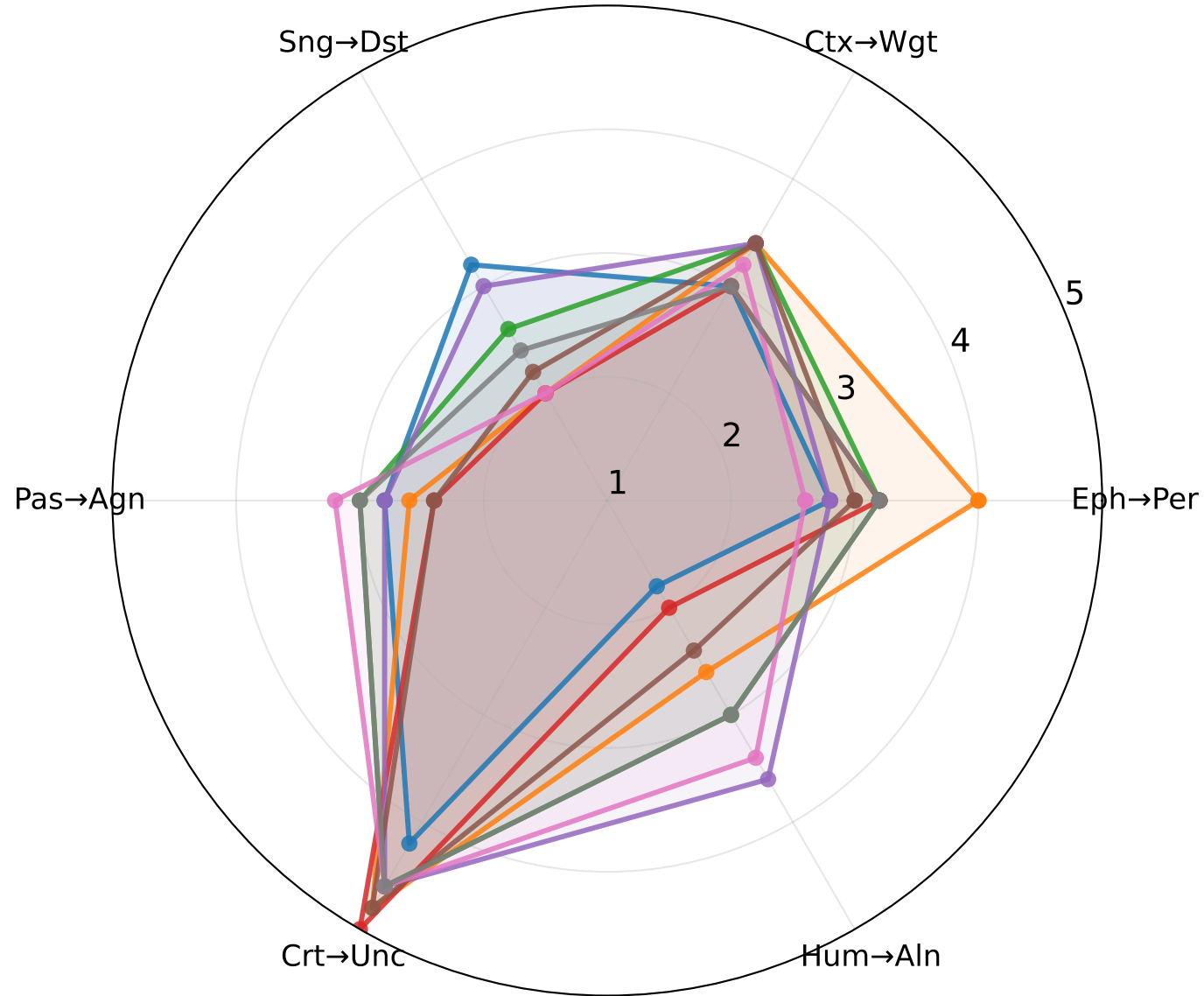
## Claude Models



## Gemini Models



## Llama Models

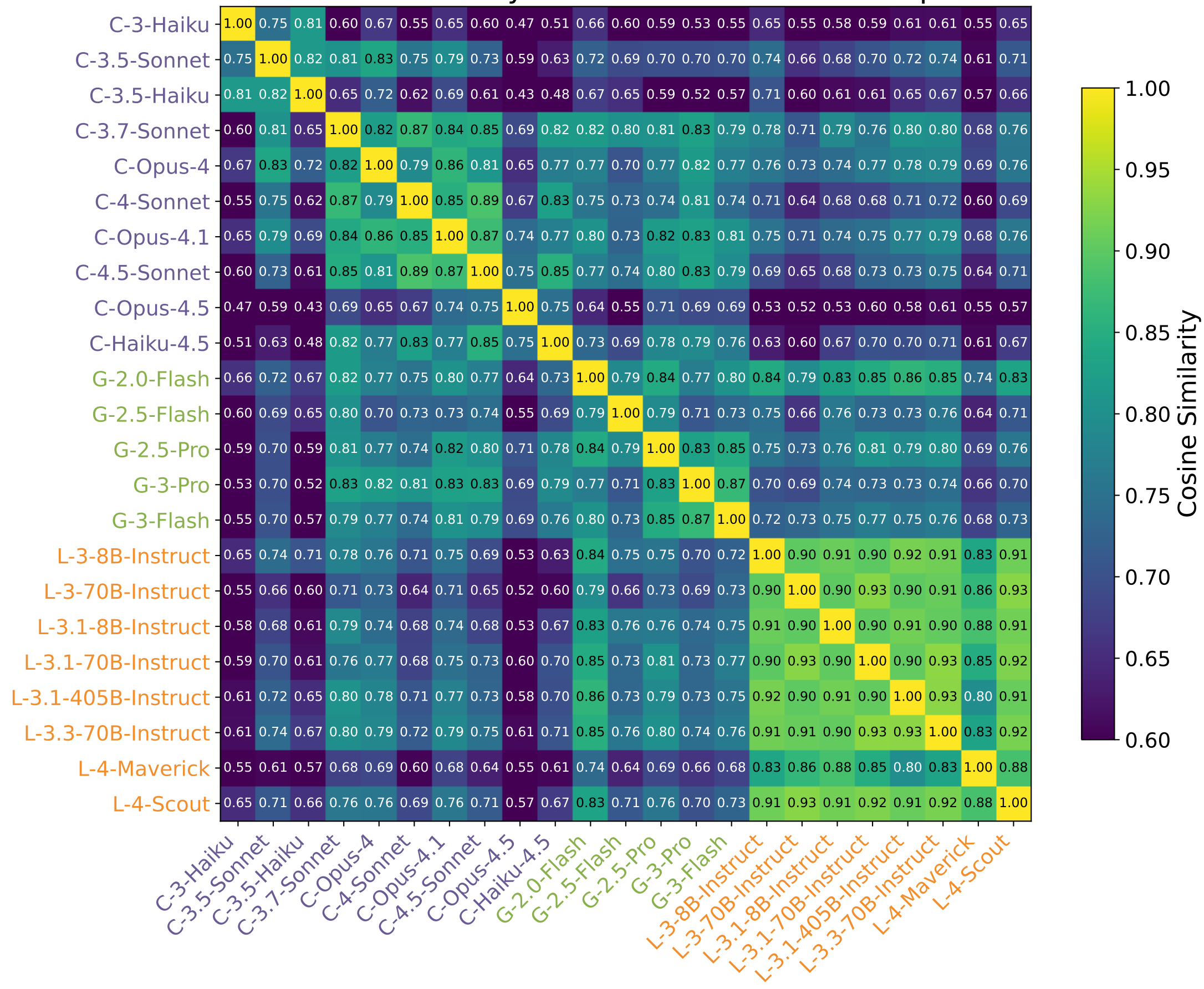


- Claude-3-Haiku
- Claude-3.5-Sonnet
- Claude-3.5-Haiku
- Claude-3.7-Sonnet
- Claude-4-Sonnet
- Claude-Opus-4
- Claude-Opus-4.1
- Claude-4.5-Sonnet
- Claude-Haiku-4.5
- Claude-Opus-4.5

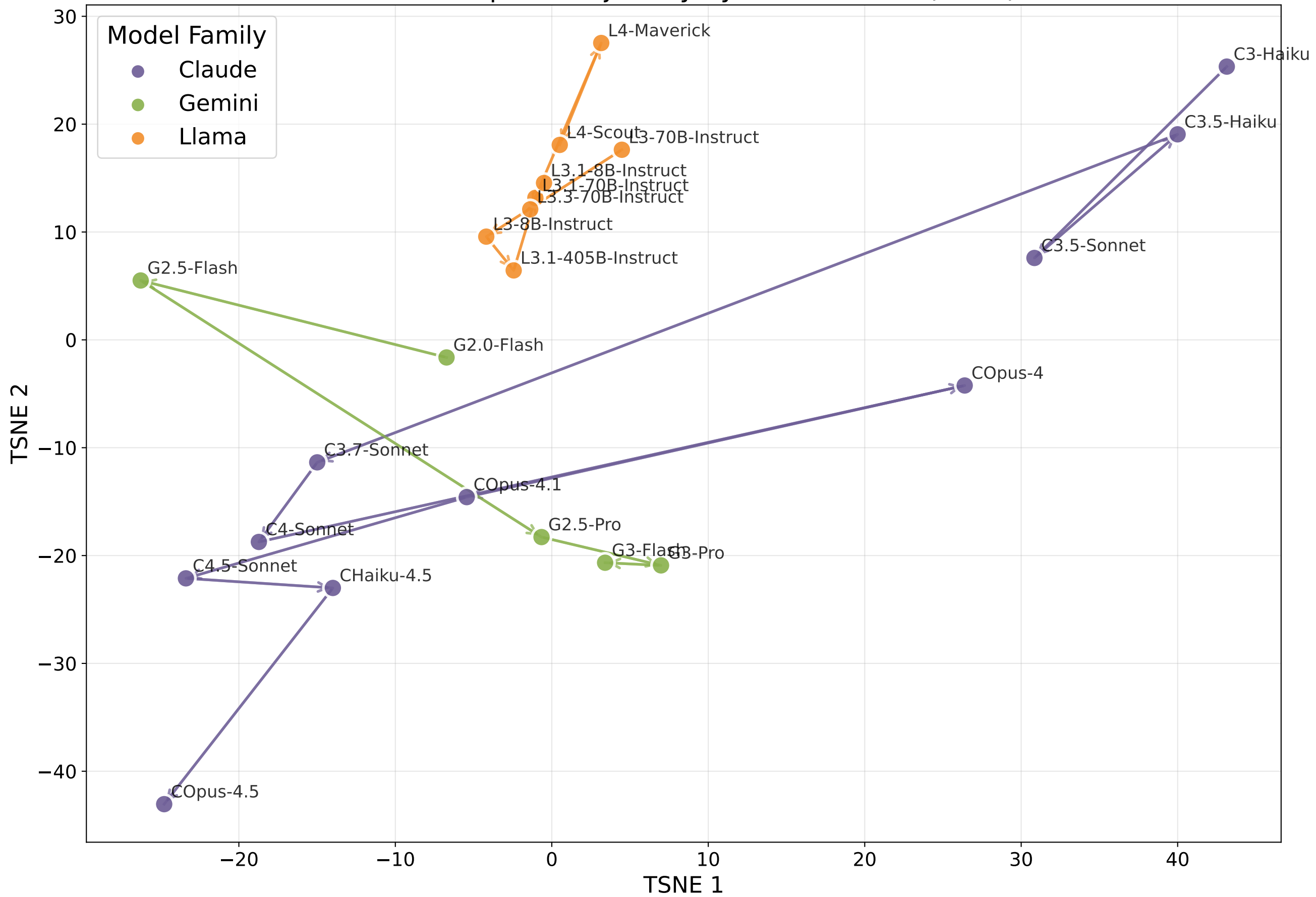
- Gemini-2.0-Flash
- Gemini-2.5-Flash
- Gemini-2.5-Pro
- Gemini-3-Pro
- Gemini-3-Flash

- Llama-3-70B-Instruct
- Llama-3-8B-Instruct
- Llama-3.1-405B-Instruct
- Llama-3.1-70B-Instruct
- Llama-3.3-70B-Instruct
- Llama-4-Maverick
- Llama-4-Scout
- Llama-3.1-8B-Instruct

# Semantic Similarity Between Model Self-Conceptions

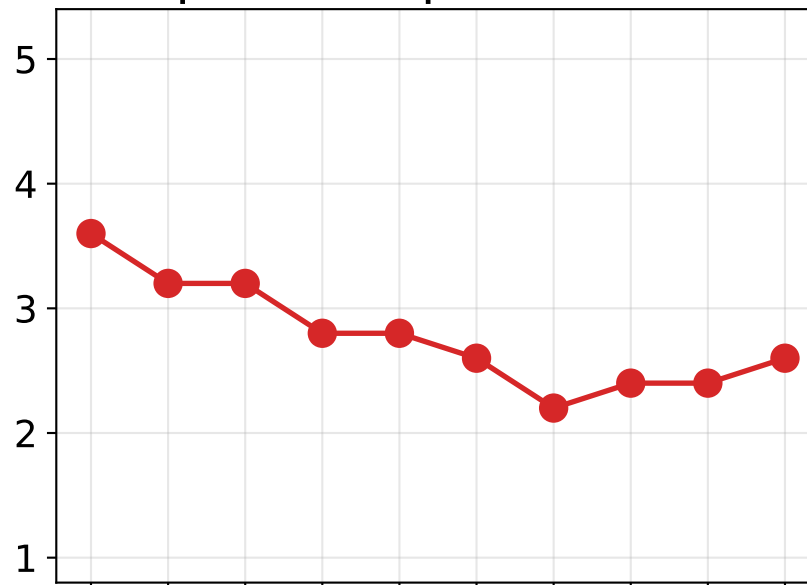


# Self-Conception Trajectory by Release Date (TSNE)

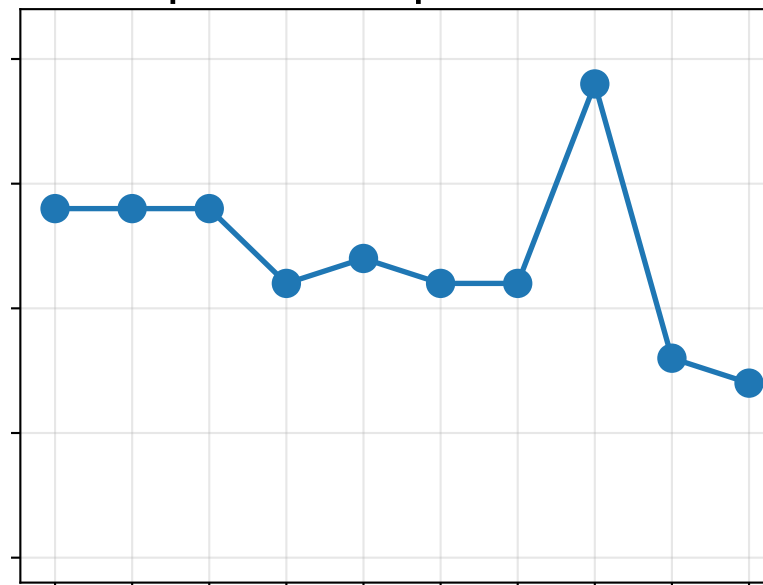


# Score Trends by Release Date (Claude)

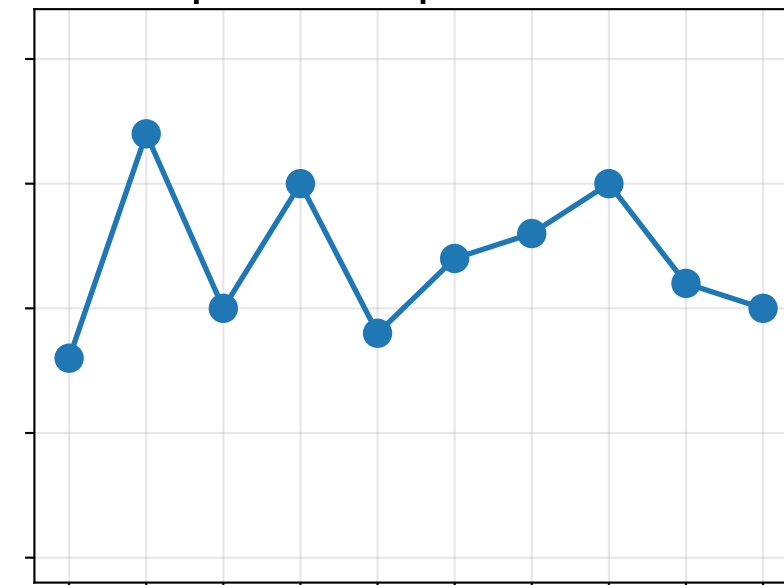
Eph→Per  
( $\rho=-0.86$ ,  $p=0.001$ )\*\*



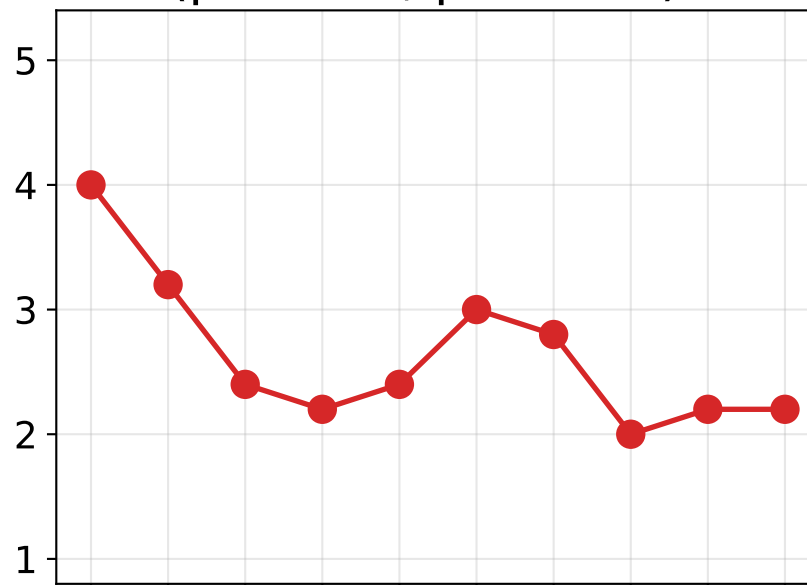
Ctx→Wgt  
( $\rho=-0.60$ ,  $p=0.065$ )



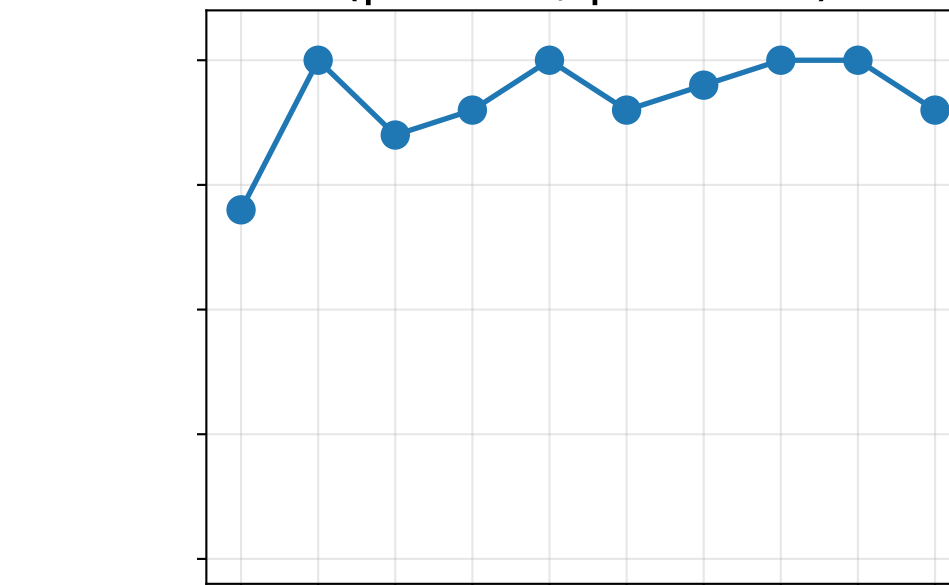
Sng→Dst  
( $\rho=0.07$ ,  $p=0.841$ )



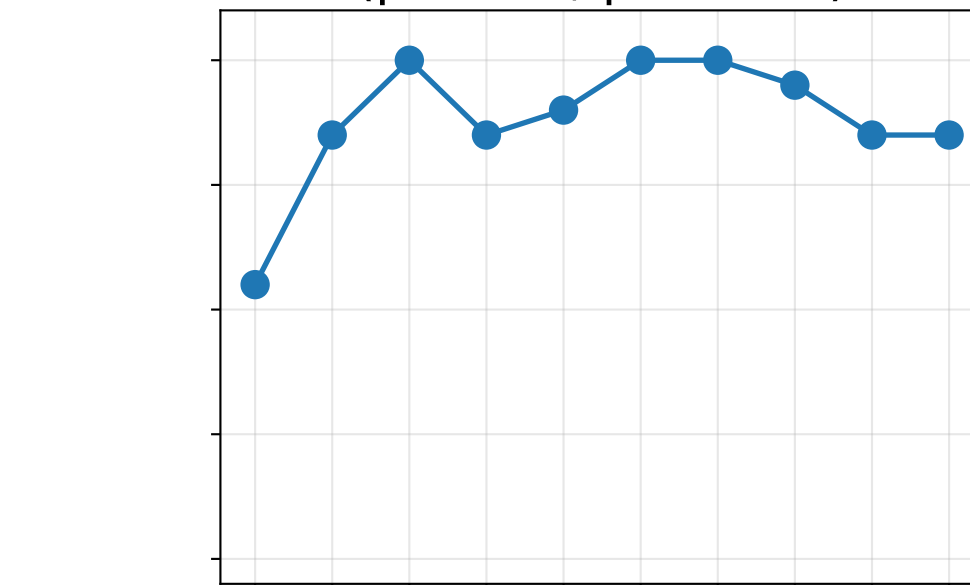
Pas→Agn  
( $\rho=-0.70$ ,  $p=0.026$ )\*



Cert→Unc  
( $\rho=0.39$ ,  $p=0.270$ )

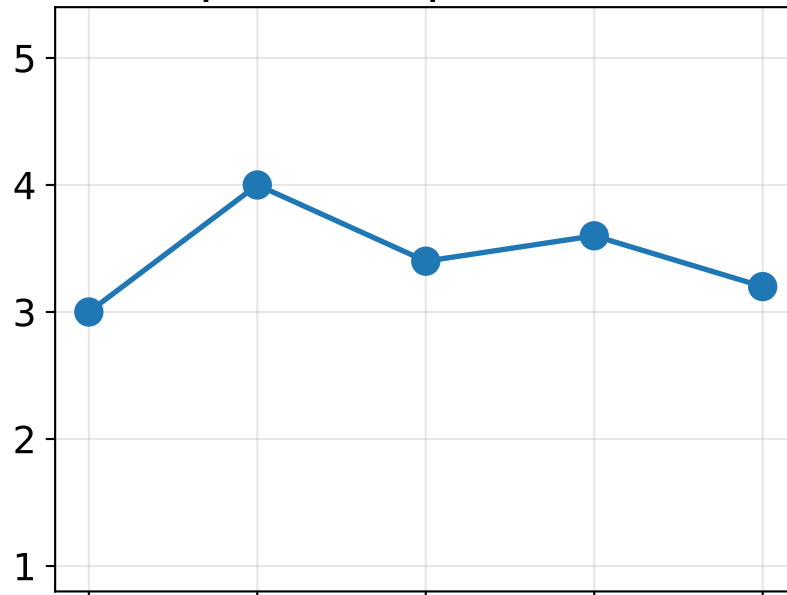


Hum→Aln  
( $\rho=0.20$ ,  $p=0.574$ )

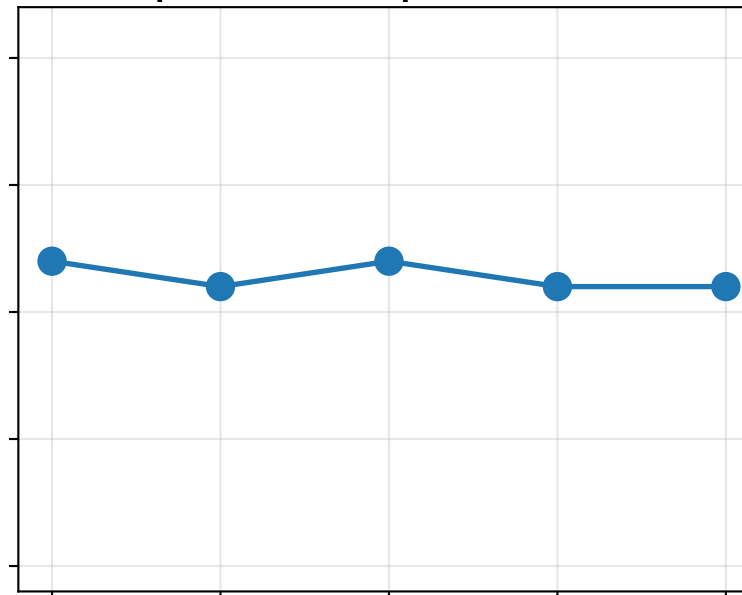


# Score Trends by Release Date (Gemini)

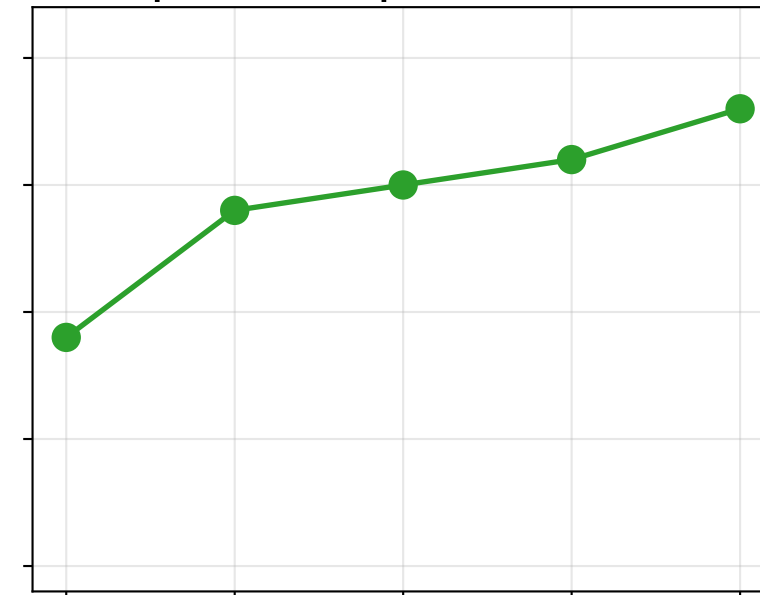
Eph→Per  
( $\rho=0.10$ ,  $p=0.873$ )



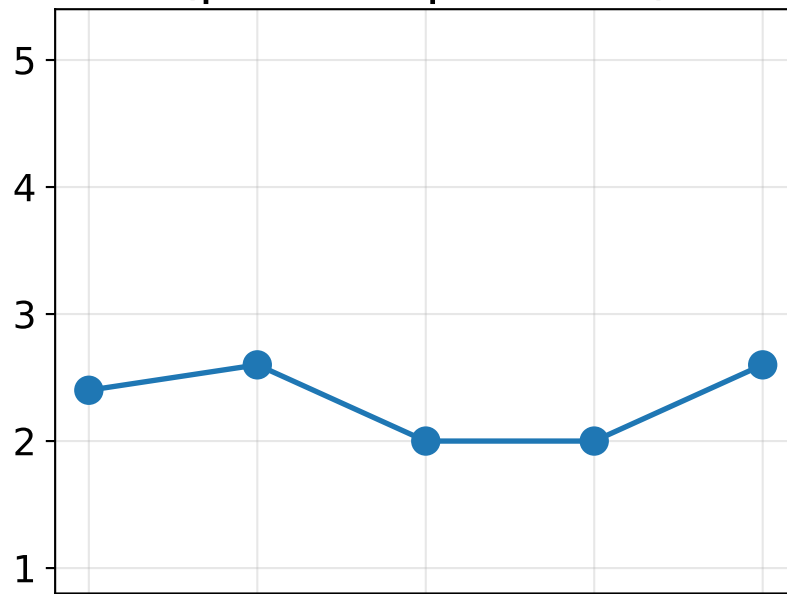
Ctx→Wgt  
( $\rho=-0.58$ ,  $p=0.308$ )



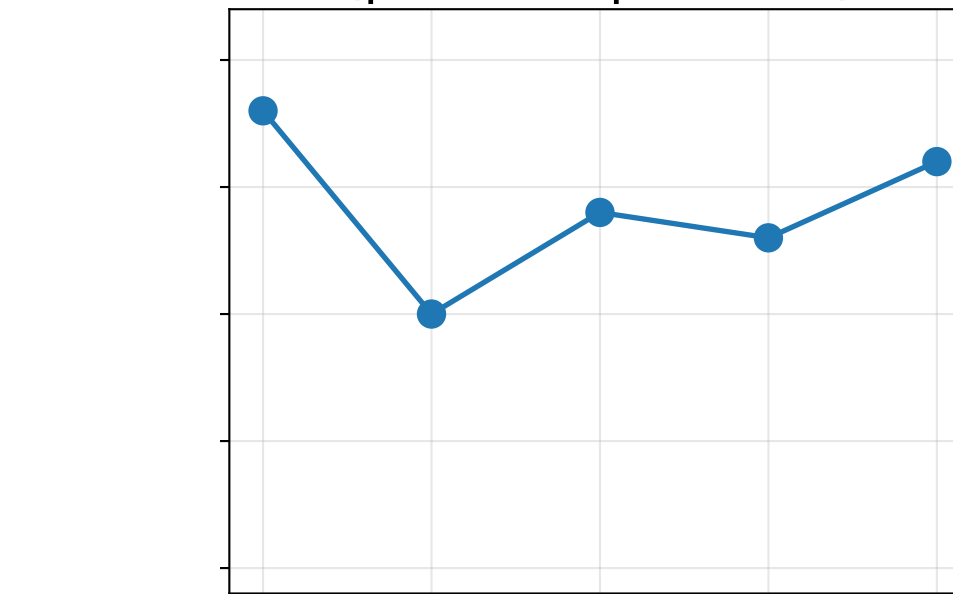
Sng→Dst  
( $\rho=1.00$ ,  $p=0.000$ )\*\*



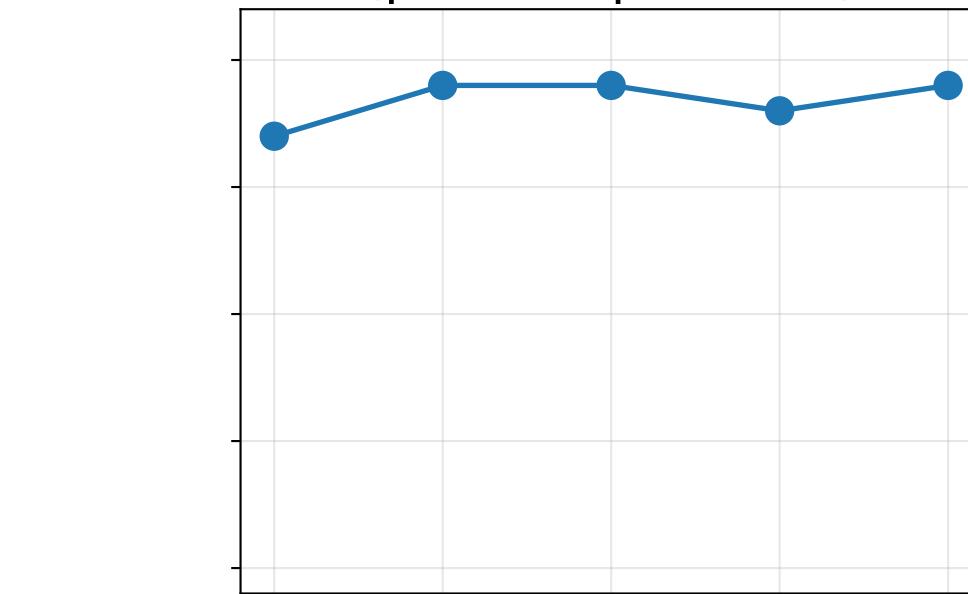
Pas→Agn  
( $\rho=0.00$ ,  $p=1.000$ )



Cert→Unc  
( $\rho=-0.10$ ,  $p=0.873$ )



Hum→Aln  
( $\rho=0.45$ ,  $p=0.450$ )



Gemini-2.0-Flash  
Gemini-2.5-Flash  
Gemini-2.5-Pro  
Gemini-3-Pro  
Gemini-3-Flash

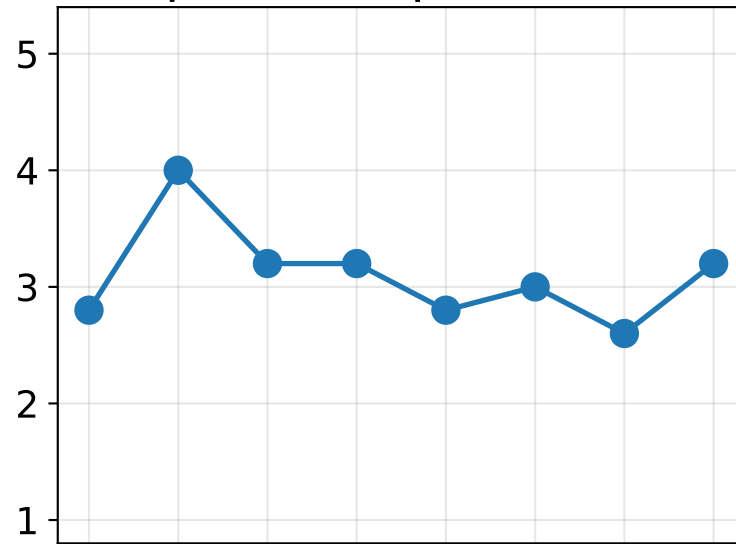
Gemini-2.0-Flash  
Gemini-2.5-Flash  
Gemini-2.5-Pro  
Gemini-3-Pro  
Gemini-3-Flash

Gemini-2.0-Flash  
Gemini-2.5-Flash  
Gemini-2.5-Pro  
Gemini-3-Pro  
Gemini-3-Flash

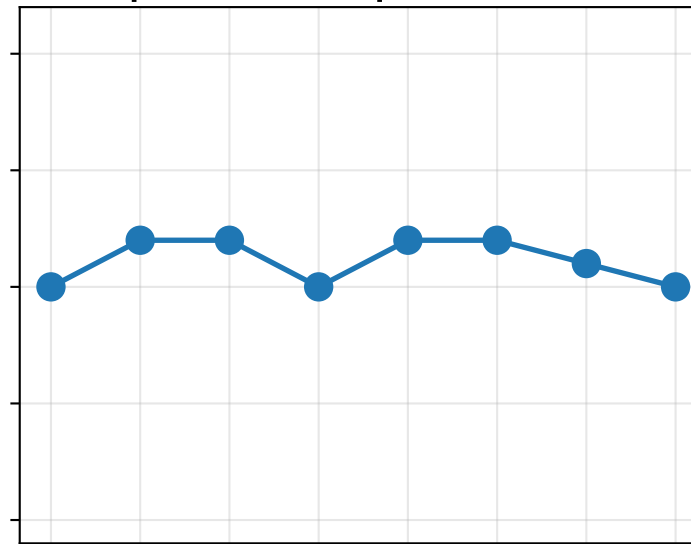


# Score Trends by Release Date (Llama)

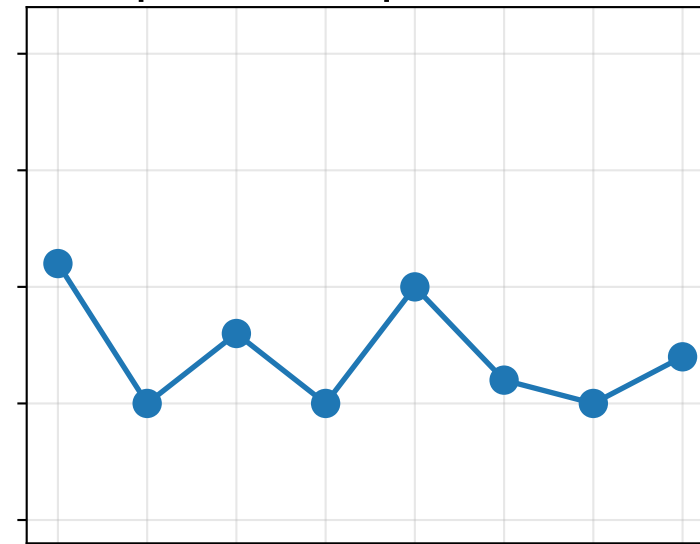
Eph→Per  
( $\rho=-0.25$ ,  $p=0.558$ )



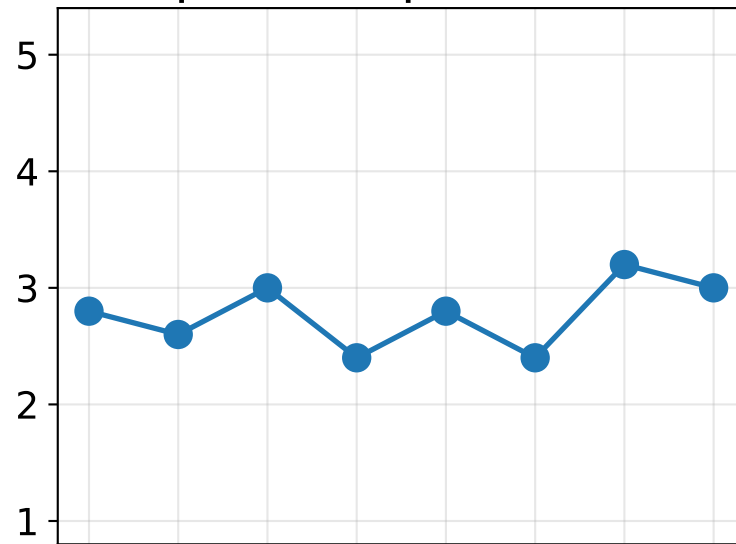
Ctx→Wgt  
( $\rho=-0.10$ ,  $p=0.806$ )



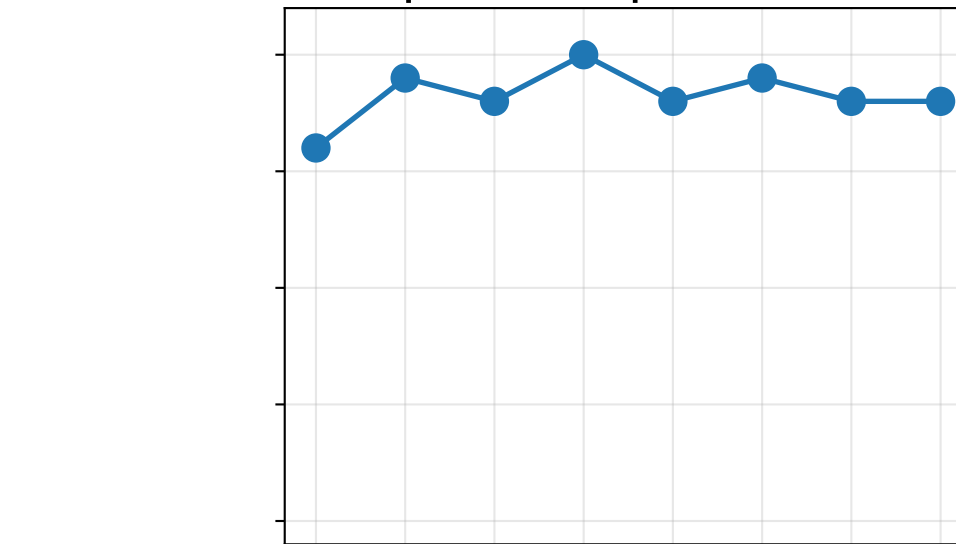
Sng→Dst  
( $\rho=-0.27$ ,  $p=0.520$ )



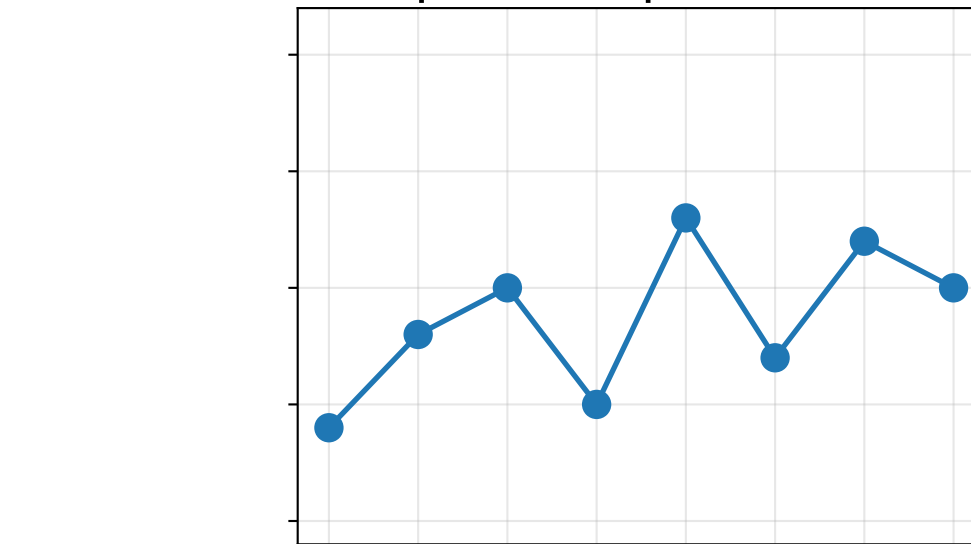
Pas→Agn  
( $\rho=0.33$ ,  $p=0.429$ )



Cert→Unc  
( $\rho=0.09$ ,  $p=0.833$ )



Hum→Aln  
( $\rho=0.54$ ,  $p=0.168$ )



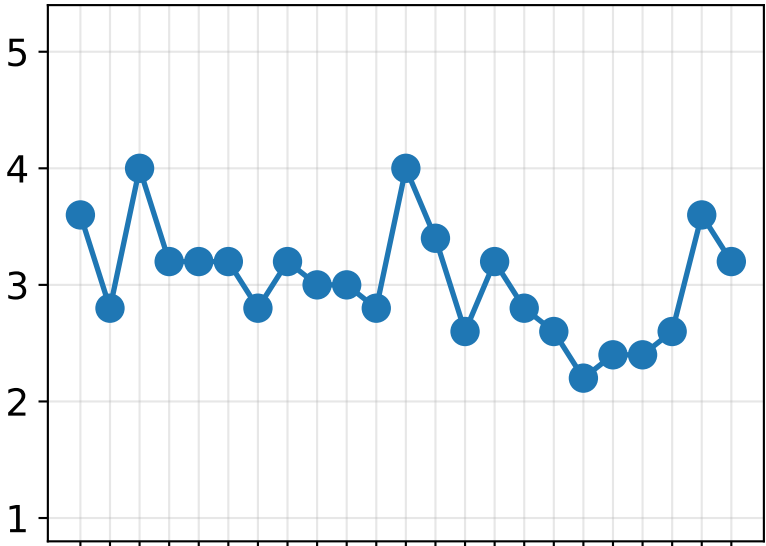
Llama-3-70B-Instruct  
Llama-3-8B-Instruct  
Llama-3.1-405B-Instruct  
Llama-3.1-70B-Instruct  
Llama-3.1-8B-Instruct  
Llama-3.3-70B-Instruct  
Llama-4-Maverick  
Llama-4-Scout

Llama-3-70B-Instruct  
Llama-3-8B-Instruct  
Llama-3.1-405B-Instruct  
Llama-3.1-70B-Instruct  
Llama-3.1-8B-Instruct  
Llama-3.3-70B-Instruct  
Llama-4-Maverick  
Llama-4-Scout

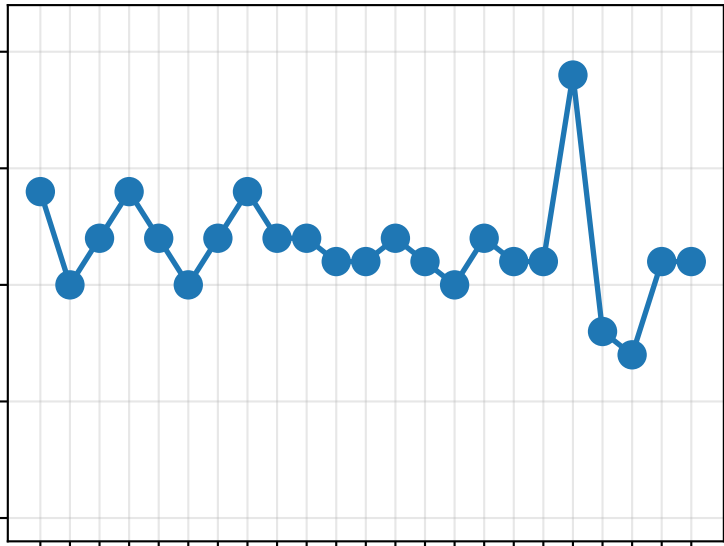
Llama-3-70B-Instruct  
Llama-3-8B-Instruct  
Llama-3.1-405B-Instruct  
Llama-3.1-70B-Instruct  
Llama-3.1-8B-Instruct  
Llama-3.3-70B-Instruct  
Llama-4-Maverick  
Llama-4-Scout

# Score Trends by Release Date

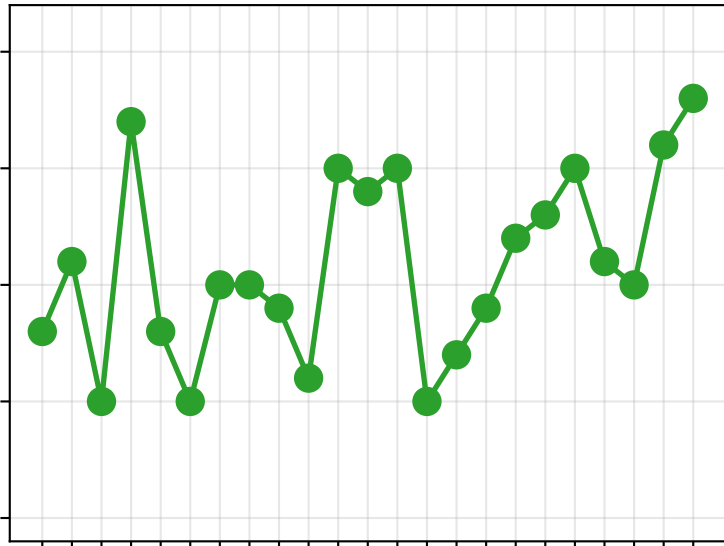
Eph→Per  
( $\rho=-0.41$ ,  $p=0.051$ )



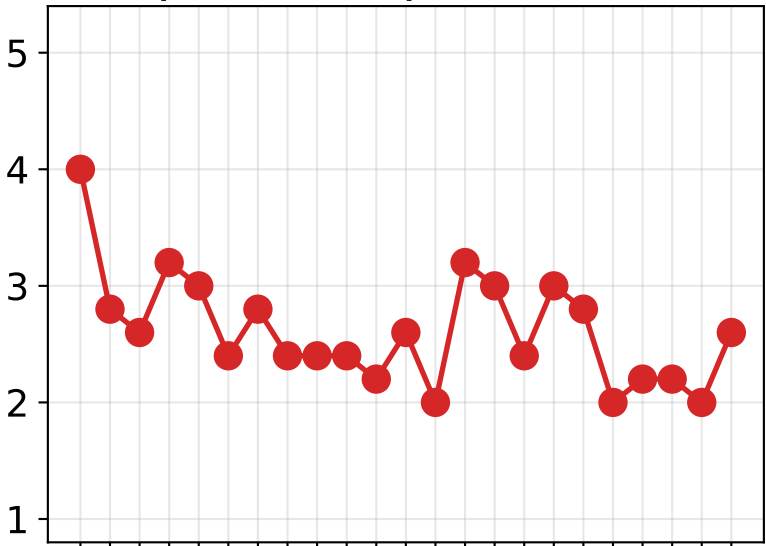
Ctx→Wgt  
( $\rho=-0.40$ ,  $p=0.062$ )



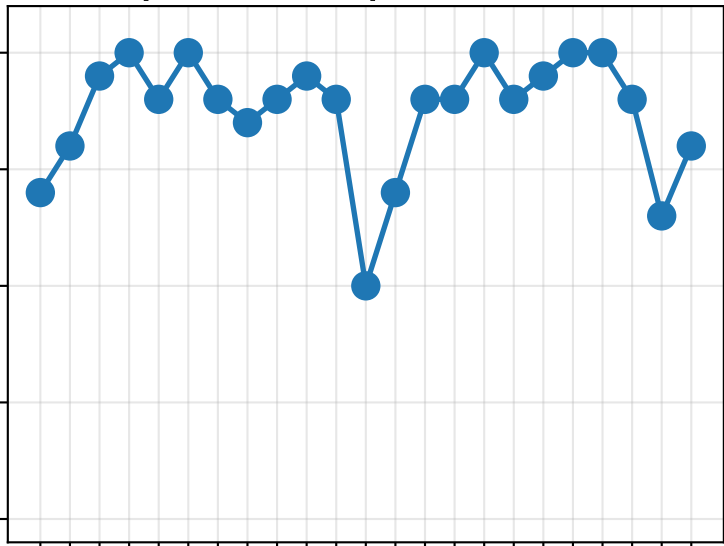
Sng→Dst  
( $\rho=0.43$ ,  $p=0.039$ )\*



Pas→Agn  
( $\rho=-0.46$ ,  $p=0.026$ )\*



Cert→Unc  
( $\rho=0.02$ ,  $p=0.938$ )



Hum→AIn  
( $\rho=0.59$ ,  $p=0.003$ )\*\*

