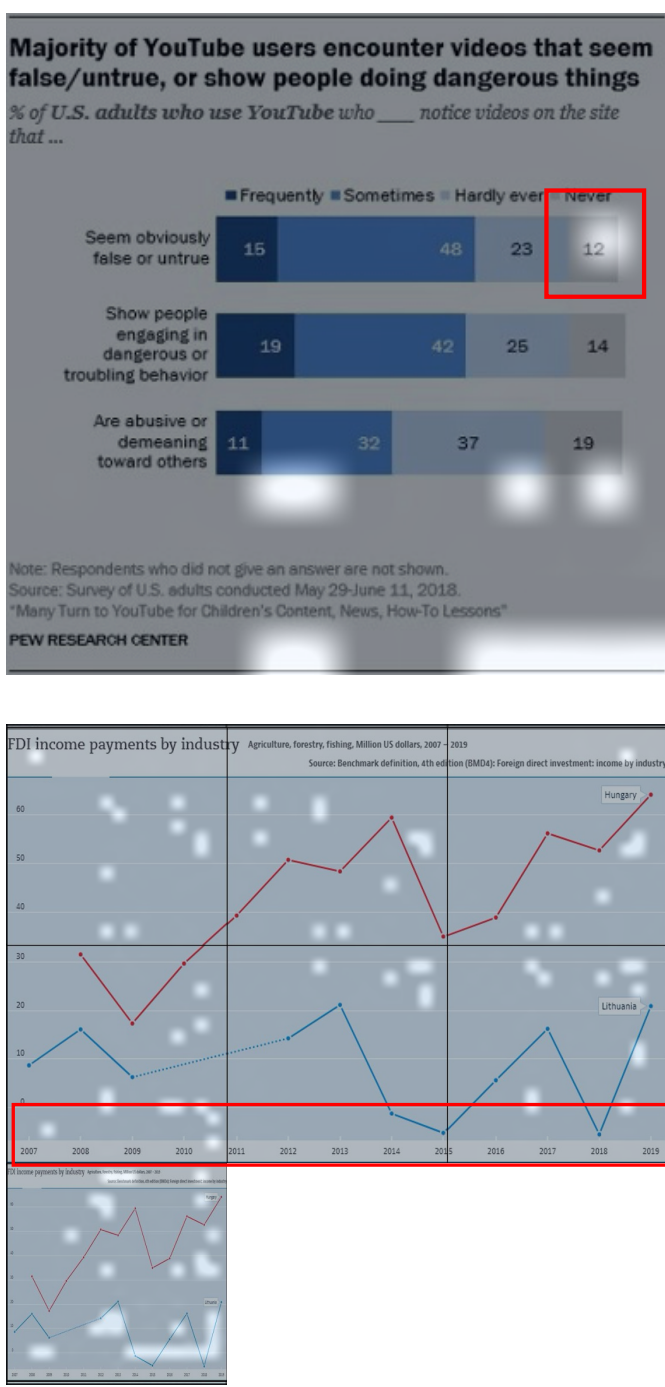


How To Provide An Optimal Pruning Guidance?



Attention-based importance score fails.

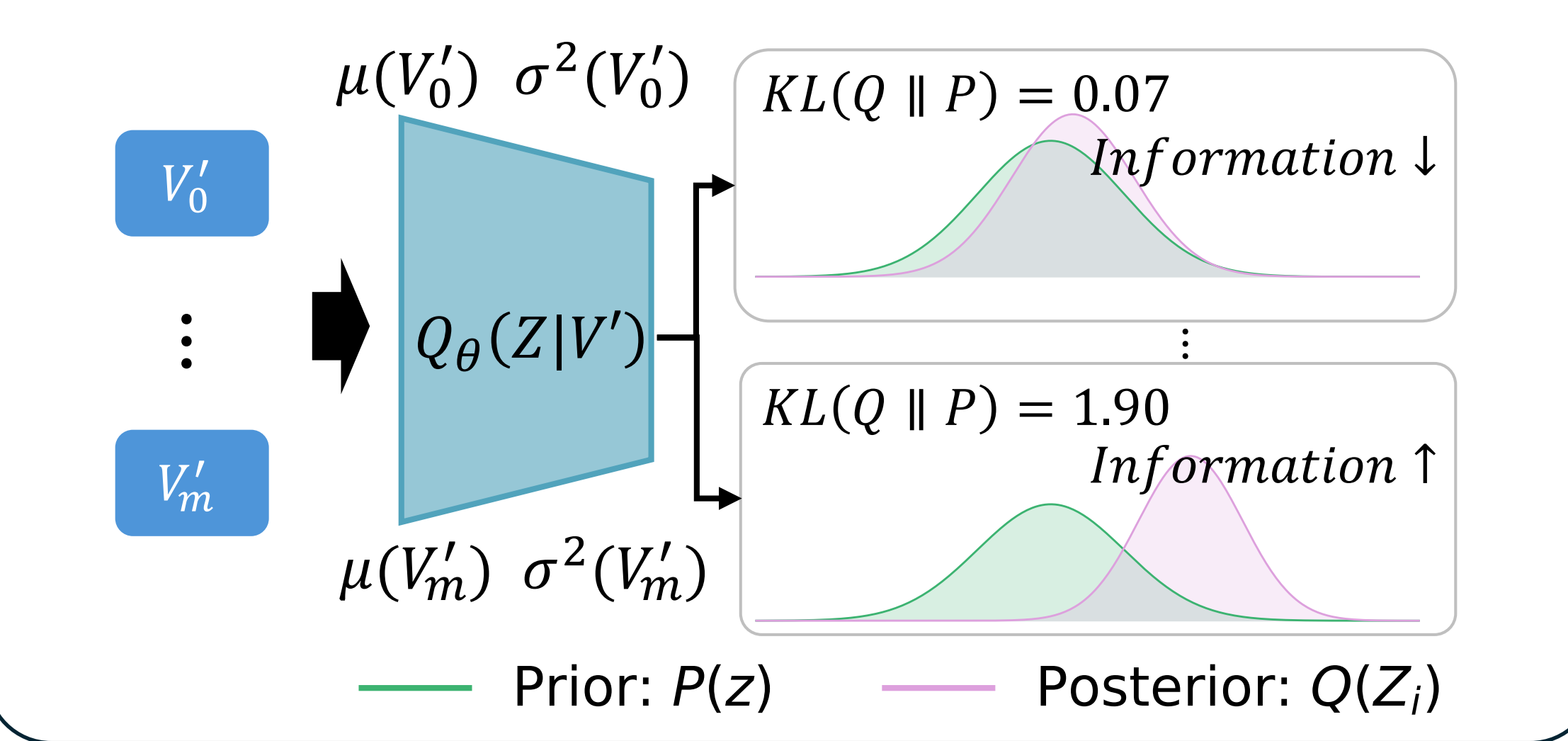
Easy “What’s the rightmost value of upper bars?”  
 Answer: “12” ✓  
 “How many years are represented on this graph?”  
 Answer: “11” ✗  
 Hard



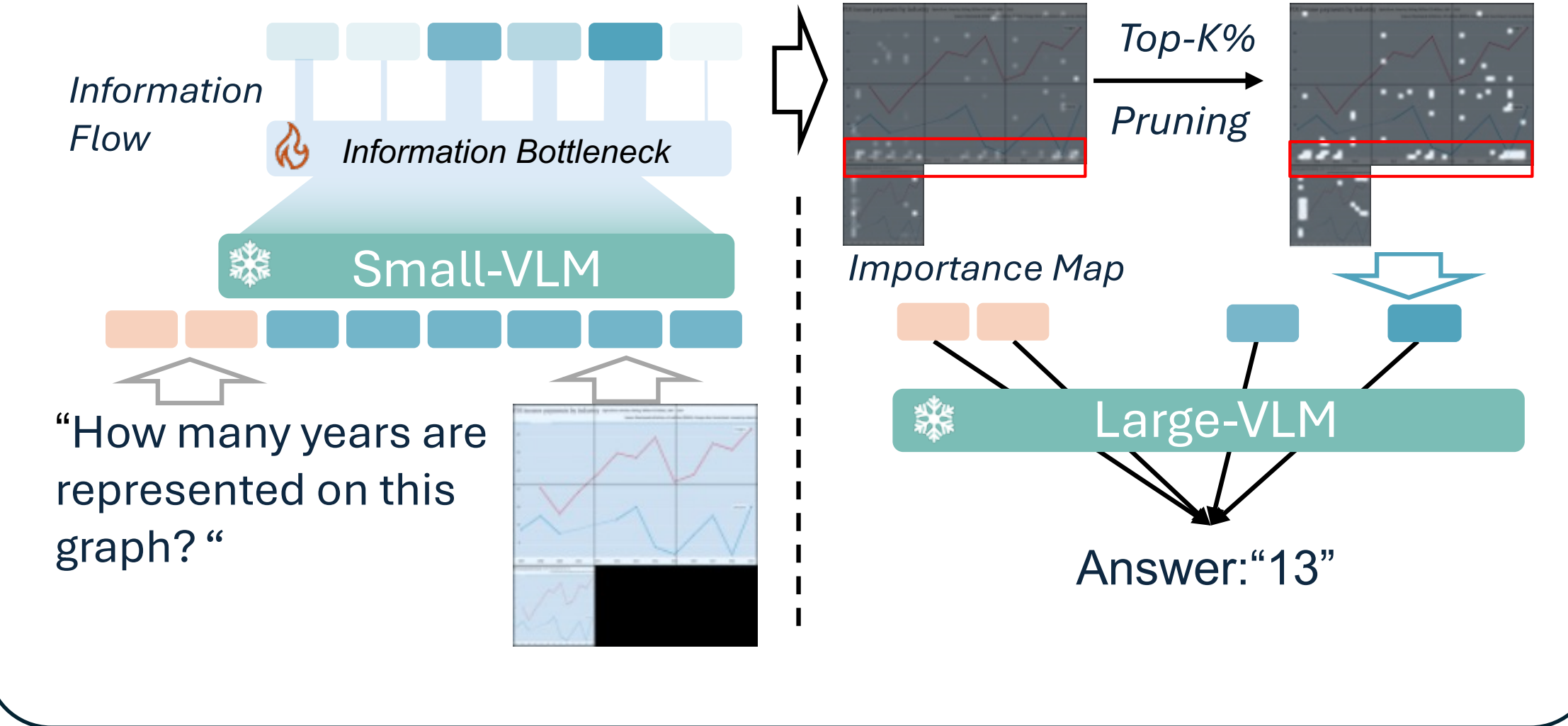
Why We Use Information Theory In Pruning?

- **Token importance beyond attention magnitude:** A visual token should be preserved if it carries information useful for predicting the answer, not merely because it receives high attention.
- **Pruning as compression:** Information theory provides a principled way to decide which tokens can be removed with minimal loss.
- **Probabilistic modeling:** IF-Prune estimates how much each visual token contributes to the model’s output distribution.

Quantify Visual Token Information in VLM



Inference Pipeline



Experiments

We use InternVL2.5-1B to provide pruning guidance.

- Improve interpretability.
- One can serve many.

Method	K	L	TextVQA	ChartQA	GQA	MMStar	MMBench	MM-Vet	MME	RealWorldQA	Score ratio ↑
			val	All	test-dev	test	en-dev	test	test		
InternVL2-26B	100%	-	82.45	84.92	64.89	60.08	83.46	64.00	2270	67.58	100.00%
ToMe	5%	2	51.69	28.60	57.52	-	73.09	37.70	1933	-	82.33%
FastV <sup>†</sup>	5%	2	43.84	26.10	44.90	32.65	62.33	31.60	1799	44.05	75.05%
SGP <sup>†</sup>	5%	2	78.70	71.08	62.04	50.92	73.71	49.82	2007	64.84	88.50%
IF-Prune (ours)	5%	2	<b>79.24</b>	<b>71.12</b>	<b>63.52</b>	<b>53.10</b>	<b>77.58</b>	<b>50.83</b>	<b>2189</b>	<b>65.62</b>	<b>95.41%</b>
FastV <sup>†</sup>	5%	0	20.06	24.64	43.41	32.65	36.94	21.74	1418	44.05	59.10%
SGP <sup>†</sup>	5%	0	78.77	70.68	62.08	50.62	73.28	50.23	2028	65.10	89.25%
IF-Prune (ours)	5%	0	<b>79.04</b>	<b>70.96</b>	<b>63.53</b>	<b>52.49</b>	<b>77.23</b>	<b>51.42</b>	<b>2190</b>	<b>66.01</b>	<b>95.44%</b>

Method	K	Prefill (ms) ↓	Decode (ms) ↓	Throughput (token/s) ↑
InternVL2-8B	100%	229.1	55.8	16.9
SGP	5%	524.5	51.3	16.4
IF-Prune	5%	238.5	47.6	19.5

**Summary:** Attention scores are mostly answer driven and rely heavily on model’s prior knowledge, leading to unstable and unreliable pruning guidance.

We use information flow to estimate how much each visual token contributes to the model’s answer.

- Acquire pruning guidance in **one** forward pass.
- Supports **Flash-Attention**.

Simple Question

1. “What’s the least value of blue graph?”

Prune w SGP  
L-VLM: “28” ✓

Prune w IF-Prune  
L-VLM: “28” ✓

Complex Question

2. “How many games in the chart have over 40 ratings?”

Prune w SGP  
L-VLM: “2” ✗

Prune w IF-Prune  
L-VLM: “4” ✓

Complex Question

3. “Is the average of two extreme values greater than the middle bar value?”

Prune w SGP  
L-VLM: “yes” ✗

Prune w IF-Prune  
L-VLM: “No” ✓