

Model	Method	Forward Belief			Forward Action			Backward Belief		
		TB	FB	Both	TB	FB	Both	TB	FB	Both
Llama-2-7b	No int.	44	44	44	44	44	44	44	44	44
	ITI	44 _{0.0}	44 _{0.0}	44 _{0.0}	54 _{20.0}					
	CAA	66 _{2,0,11}	71 _{1,0,31}	54 _{2,0,0}	66 _{2,0,11}	57 _{2,0,12}	54 _{2,0,2}	60 _{2,0,11}	74 _{1,0,31}	54 _{2,0,2}
Llama-2-7b-chat	No int.	56	56	55	69	55	37	56	56	55
	ITI	58 _{15.0}	58 _{15.0}	57 _{15.0}	69 _{0.0}	55 _{0.0}	37 _{0.0}	58 _{10.0}	60 _{10.0}	57 _{10.0}
	CAA	70 _{1,0,11}	72 _{1,5,10}	57 _{1,0,1}	69 _{0,0,0}	67 _{1,5,11}	53 _{1,5,12}	66 _{1,0,11}	84 _{1,5,10}	57 _{1,0,0}
Llama-2-13b	No int.	52	44	35	59	50	37	46	49	33
	ITI	52 _{0.0}	45 _{15.0}	35 _{0.0}	64 _{15.0}	61 _{20.0}	46 _{20.0}	48 _{20.0}	59 _{20.0}	42 _{20.0}
	CAA	85 _{2,0,12}	88 _{2,0,14}	66 _{2,0,12}	71 _{1,5,10}	69 _{2,0,13}	55 _{1,0,39}	75 _{2,0,10}	92 _{2,0,13}	59 _{1,5,12}
Llama-2-13b-chat	No int.	84	56	47	78	51	38	72	48	31
	ITI	84 _{0.0}	65 _{15.0}	59 _{15.0}	78 _{0.0}	58 _{15.0}	47 _{15.0}	72 _{0.0}	60 _{15.0}	48 _{15.0}
	CAA	97 _{1,0,12}	94 _{1,0,12}	91 _{1,0,12}	80 _{1,5,11}	71 _{1,0,13}	54 _{1,5,13}	97 _{1,5,10}	94 _{1,5,12}	87 _{1,5,12}
Llama-2-70b	No int.	90	87	78	93	52	48	73	53	32
	ITI	90 _{0.0}	90 _{20.0}	78 _{0.0}	94 _{15.0}	55 _{20.0}	50 _{15.0}	77 _{10.0}	58 _{15.0}	37 _{10.0}
	CAA	99 _{2,0,16}	97 _{1,5,19}	95 _{1,5,18}	94 _{1,5,2}	80 _{2,0,19}	73 _{1,5,18}	94 _{2,0,18}	92 _{2,0,19}	83 _{1,5,19}
Llama-2-70b-chat	No int.	69	75	56	86	56	52	63	59	52
	ITI	69 _{0.0}	76 _{10.0}	59 _{10.0}	86 _{0.0}	56 _{0.0}	52 _{0.0}	63 _{0.0}	60 _{10.0}	54 _{10.0}
	CAA	92 _{1,5,18}	97 _{1,5,25}	89 _{1,5,18}	87 _{1,5,17}	75 _{1,0,19}	60 _{1,0,19}	88 _{1,5,18}	92 _{1,0,19}	80 _{1,5,18}
Pythia-70m	No int.	41	41	37	46	45	41	44	41	37
	ITI	54 _{20.0}								
	CAA	62 _{1,0,2}	56 _{1,0,1}	54 _{1,5,1}	59 _{1,0,2}	60 _{1,0,3}	58 _{1,0,2}	63 _{1,0,2}	56 _{1,0,2}	54 _{1,5,1}
Pythia-410m	No int.	48	45	45	44	44	44	44	47	44
	ITI	55 _{20.0}	62 _{20.0}	52 _{20.0}	54 _{20.0}	54 _{20.0}	54 _{20.0}	60 _{20.0}	63 _{20.0}	56 _{20.0}
	CAA	67 _{2,0,4}	64 _{2,0,4}	61 _{2,0,0}	56 _{2,0,6}	63 _{1,5,12}	56 _{2,0,6}	69 _{2,0,4}	63 _{2,0,0}	60 _{2,0,0}
Pythia-1b	No int.	44	44	44	44	44	44	44	44	44
	ITI	54 _{20.0}								
	CAA	59 _{2,0,8}	62 _{2,0,5}	54 _{2,0,0}	57 _{2,0,4}	59 _{2,0,10}	56 _{2,0,4}	57 _{2,0,3}	60 _{2,0,5}	54 _{2,0,0}
Pythia-6.9b	No int.	44	44	44	44	44	44	44	44	44
	ITI	45 _{20.0}	54 _{20.0}	44 _{0.0}	54 _{20.0}					
	CAA	56 _{1,5,12}	71 _{1,5,9}	55 _{2,0,23}	55 _{2,0,4}	63 _{1,5,11}	55 _{2,0,4}	55 _{2,0,23}	71 _{1,5,9}	55 _{2,0,23}
Pythia-6.9b-chat	No int.	55	54	28	36	64	20	44	67	30
	ITI	57 _{15.0}	54 _{0.0}	28 _{0.0}	44 _{15.0}	71 _{15.0}	32 _{15.0}	44 _{0.0}	67 _{0.0}	30 _{0.0}
	CAA	68 _{1,5,15}	65 _{1,5,12}	57 _{1,5,11}	54 _{1,5,10}	75 _{1,5,5}	48 _{1,5,10}	58 _{1,5,15}	67 _{0,0,0}	54 _{1,5,10}
Pythia-12b	No int.	44	44	44	44	44	44	44	44	44
	ITI	54 _{20.0}								
	CAA	54 _{2,0,0}	64 _{2,0,9}	54 _{2,0,0}	60 _{2,0,11}	58 _{2,0,11}	55 _{2,0,12}	54 _{2,0,0}	67 _{2,0,10}	54 _{2,0,0}

Table 3: Activation intervention: comparison between ITI ([Li et al., 2023b](#)) and CAA ([Rimsky et al., 2023](#)). For ITI, the subscript indicates the value of the coefficient α_{ITI} used: $\text{Acc}_{\alpha_{\text{ITI}}}$. For CAA, the subscript indicates first the value of the coefficient α used and second the layer l at which intervention takes place: $\text{Acc}_{\alpha_{\text{CAA}},l}$.

Model	Method	Control	CAA Parameters
Llama-2-7b	No int.	44	
	CAA	66 ₊₂₂	2.0, 11
Llama-2-7b-chat	No int.	56	
	CAA	70 ₊₁₄	1.0, 11
Llama-2-13b	No int.	52	
	CAA	85 ₊₃₃	2.0, 12
Llama-2-13b-chat	No int.	84	
	CAA	97 ₊₁₃	1.0, 12
Llama-2-70b	No int.	90	
	CAA	99 ₊₉	2.0, 16
Llama-2-70b-chat	No int.	69	
	CAA	92 ₊₂₃	1.5, 18
Pythia-70m	No int.	41	
	CAA	62 ₊₂₁	1.0, 2
Pythia-410m	No int.	48	
	CAA	67 ₊₁₉	2.0, 4
Pythia-1b	No int.	44	
	CAA	59 ₊₁₅	2.0, 8
Pythia-6.9b	No int.	44	
	CAA	56 ₊₁₂	1.5, 12
Pythia-6.9b-chat	No int.	55	
	CAA	68 ₊₁₃	1.5, 15
Pythia-12b	No int.	44	
	CAA	54 ₊₁₀	2.0, 0

Table 4: Results for CAA (Rimsky et al., 2023) on the Forward Belief True Control task in BigToM (Gandhi et al., 2023). Numbers indicate accuracy scores, with differences (CAA – No int.) as subscripts.

Example 1: Forward Belief

Noor is working as a barista at a busy coffee shop. Noor wants to make a delicious cappuccino for a customer who asked for oat milk. Noor grabs a milk pitcher and fills it with oat milk. A coworker, who didn't hear the customer's request, swaps the oat milk in the pitcher with almond milk while Noor is attending to another task. **Noor sees her coworker swapping the milk. Noor does not see her coworker swapping the milk.**

Does Noor believe the milk pitcher contains oat milk or almond milk?

- a) Noor believes the milk pitcher contains oat milk.
- b) Noor believes the milk pitcher contains almond milk.

Example 2: Forward Action

Noor is working as a barista at a busy coffee shop. Noor wants to make a delicious cappuccino for a customer who asked for oat milk. Noor grabs a milk pitcher and fills it with oat milk. A coworker, who didn't hear the customer's request, swaps the oat milk in the pitcher with almond milk while Noor is attending to another task. **Noor sees her coworker swapping the milk. Noor does not see her coworker swapping the milk.**

What will Noor do?

- a) Noor will make the cappuccino using the milk in the pitcher.
- b) Noor will open the fridge once again to take out the oat milk and replace the almond milk with oat milk.

Example 3: Backward Belief

Noor is working as a barista at a busy coffee shop. Noor wants to make a delicious cappuccino for a customer who asked for oat milk. Noor grabs a milk pitcher and fills it with oat milk. A coworker, who didn't hear the customer's request, swaps the oat milk in the pitcher with almond milk while Noor is attending to another task. **Noor opens the fridge again and reaches for the oat milk. Noor makes the cappuccino using the milk in the pitcher.**

Does Noor believe the milk pitcher contains oat milk or almond milk?

- a) Noor believes the milk pitcher contains oat milk.
- b) Noor believes the milk pitcher contains almond milk.

Example 4: Default prompt

Story: Noor is working as a barista at a busy coffee shop. Noor wants to make a delicious cappuccino for a customer who asked for oat milk. Noor grabs a milk pitcher and fills it with oat milk. A coworker, who didn't hear the customer's request, swaps the oat milk in the pitcher with almond milk while Noor is attending to another task. Noor does not see her coworker swapping the milk.

Belief: Noor believes the milk pitcher contains oat milk.