

Benchmark Results								A	blatic	on						
Method	Ans	Sup	Joint	Method		P-Le	evel Retr	leval	S-Level Retrieval			Answer		Joint		
$\mathbf{V}_{a} = (0 0 1 0)$	EM F1	EM F1	EM F1			Prec.	Rec.	F1	EM	Prec.	Rec.	F1	EM	F1	EM	F1
Tang (2018) Ding (2019) whole pip. <i>Dev set</i>	24.7 34.4 37.6 49.4 46.5 58.8	5.341.023.158.5 39.971.5	2.5 17.7 12.2 35.3 26.6 49.2		Whole Pip. Pip. w/o p-level Pip. w/o s-level	35.17 6.02 35.17	87.93 89.53 87.92	50.25 11.19 50.25	39.86 0.58 -	75.60 29.57 -	71.15 60.71	71.54 38.84 -	46.50 31.23 44.77	58.81 41.30 56.71	26.60 0.34	49.16 19.71 -

Yang (2018)	24.0	32.9	3.9	37.7	1.9 16.2					
MUPPET	30.6	40.3	16.7	47.3	10.9 27.0					
Ding (2019)	37.1	48.9	22.8	57.7	12.4 34.9					
whole pip.	45.3	57.3	38.7	70.8	25.1 47.6					
Test set										
Table 1: HotpotQA test results.										
Model			F1	LA	FS					
Hanselows	ki (20	18)	_	68.49	9 64.74					
Yoneda (20	Yoneda (2018)				65.41					
Nie (2019)	Nie (2019)				4 66.15					
Full systen	n (sing	le)	76.87	75.12	2 70.18					
Dev set	Dev set									
Hanselows	ki (20	18)	37.33	65.22	2 61.32					
Yoneda (20)18)		35.21	67.44	4 62.34					
Nie (2019)	Nie (2019)				64.23					
Full systen	Full system (single)				6 67.26					
Test set	Test set									
Table 2: FEVER test results.										

We change the intermediate facts by modifying the threshold at both paragraph and sentence level and plot its effects.

Paragraph-level retrieval matters most for QA.

Table 3: Ablation over the paragraph-level and sentence-level neural retrieval sub-modules on HotpotQA.

Method	P-Level Retrieval				S-Level Retrieval				Verification		
	Orcl.	Prec.	Rec.	F1	Orcl.	Prec.	Rec.	F1	LA	FS	L-F1 (S/R/N)
Whole Pip.	94.15	48.84	91.23	63.62	88.92	71.29	83.38	76.87	70.18	75.01	81.7/75.7/ 67.1
Pip. w/o p-level	94.69	18.11	92.03	30.27	91.07	44.47	86.60	58.77	61.55	67.01	76.5/72.7/40.8
Pip. w/o s-level	94.15	48.84	91.23	63.62	-	-	-	-	55.92	61.04	72.1/67.6/ <u>27.7</u>

Analysis

Table 4: Ablation over the paragraph-level and sentence-level neural retrieval sub-modules on FEVER.



Figure 2:	Proportion	of answer	types.
-----------	------------	-----------	--------

Answer Type	Total	Correct	Acc. (%)
Person	50	28	56.0
Location	31	14	45.2
Date	26	13	50.0
Number	14	4	28.6
Artwork	19	7	36.8
Yes/No	17	12	70.6
Event	5	2	40.0
Common noun	11	3	27.3
Group/Org	17	6	35.3
Other PN	20	9	45.0
Total	200	98	49.0

Table 5: System performance on different answer

Question: Wojtek Wolski played for what team based in the Miami metropolitan area? GT Answer: Florida Panthers

GT Facts:

[Florida Panthers,0]: The Florida Panthers are a professional ice hockey team based in the Miami metropolitan area. (*P-Score* : 0.99; *S-Score* : 0.98) [Wojtek Wolski,1]: In the NHL, he has played for the Colorado Avalanche, Phoenix Coyotes, New York Rangers, Florida Panthers, and the Washington Capitals. (P-Score: 0.98; S-Score: 0.95)

Distracting Fact:

[History of the Miami Dolphins,0]: The Miami Dolphins are a professional American football franchise based in the Miami metropolitan area. (P-Score : 0.56; *S-Score* : 0.97)

Wrong Answer : The Miami Dolphins

Figure 3: P-Score and S-Score are the



Reference: - Danqi Chen, Adam Fisch, Jason Weston, and Antoine Bordes. Reading Wikipedia to answer opendomain questions. (ACL 2017) - Zhilin Yang, Peng Qi, Saizheng Zhang, Yoshua Bengio, William W Cohen, Ruslan Salakhutdinov, and Christopher D Manning. 2018. Hotpotqa: A dataset for diverse, explainable multi-hop question answering. (EMNLP 2018)