

Many Hands Make Light Work: Task-Oriented Dialogue System with Module-Based Mixture-of-Experts

Anonymous ACL submission

Abstract

Task-oriented dialogue systems are broadly used in virtual assistants and other automated services, providing interfaces between users and machines to facilitate specific tasks. For example, in the context of hotel reservations, these systems not only recommend hotels that align with user preferences but also retain user requirements for future reference. Corresponding to a wide range of properties and applications of task-oriented dialogue systems, their outputs may also be diverse. Nowadays, task-oriented dialogue systems have benefited greatly from pre-trained language models (PLMs). While being effective and performant, scaling these models is expensive and complex. To address these challenges, we propose SMETOD to generate diverse natural language outputs, which scales the capacity of a task-oriented dialogue system while maintaining efficient inference. We extensively evaluate our model on dialogue state tracking, dialogue response generation, and intent prediction. Experimental results demonstrate that SMETOD consistently achieves state-of-the-art or comparable performance on all evaluated datasets. Furthermore, SMETOD shows an advantage in the cost of inference compared to existing approaches.

1 Introduction

Task-oriented dialogue systems play a crucial role in virtual assistants and various automated services through human-machine interactions. The fundamental objective of a task-oriented dialogue system is to aid users in completing specific services or tasks all achieved through natural language dialogues (Wen et al., 2017). Considering a broad range of applications, task-oriented dialogue systems should generate diverse types of outputs for processing information, evaluating user intentions, or retaining for future reference. In real-world scenarios, useful information processed from dialogue could be presented in various formats, in-

cluding form-based (Goddeau et al., 1996; Eric and Manning, 2017b), probability-based (Thomson and Young, 2010; Mrkšić et al., 2016; Lee et al., 2019), or text-based (Hosseini-Asl et al., 2020; Wang et al., 2022). Typically, several components are responsible for managing a variety of information: natural language understanding (NLU) for comprehending and translating user intent into either natural language or a format suitable for machine processing, dialogue state tracking (DST) for discerning the user’s requirements and providing a foundation for subsequent decisions, and natural language generation (NLG) generate a natural language response to the user based on the machine’s decision of the next move.

This leads to two predominant system designs, namely pipeline-based and end-to-end, divided by whether the machine-generated response is based on dialogue utterances or processed information from other components only. Either system design presents its own set of limitations in effectively addressing diverse output objectives (Takanobu et al., 2020). Drawbacks of pipeline-based systems lie in the potential for error propagation from one module to another, and local decisions can have adverse global effects (Su et al., 2016). End-to-end dialogue systems, on the other hand, raise concerns about missing all essential information that may be required other than responses. Moreover, diagnosing and considering component-flow characteristics can be challenging in end-to-end systems (Bang et al., 2023).

Despite the limitations in dialogue-system designs, there are also significant constraints in terms of scaling dialogue models with efficiency. Recent advancements have leveraged the transfer learning capabilities of pre-trained language models (PLMs) (Devlin et al., 2018; Dong et al., 2019; Radford et al., 2019; Raffel et al., 2020b) by fine-tuning (Budzianowski and Vulić, 2019; Hosseini-Asl et al., 2020; Heck et al., 2020) or pre-training di-

043
044
045
046
047
048
049
050
051
052
053
054
055
056
057
058
059
060
061
062
063
064
065
066
067
068
069
070
071
072
073
074
075
076
077
078
079
080
081
082
083

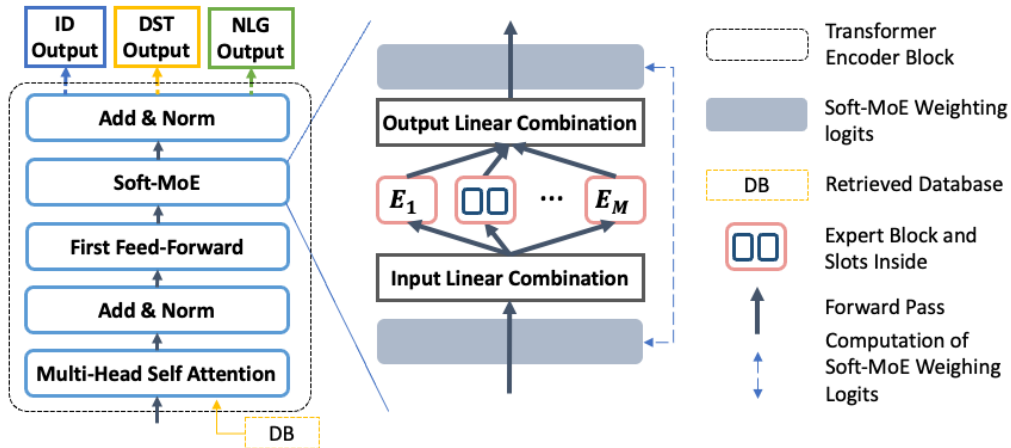


Figure 1: Architecture of the **SMETOD** as in Transformer (Vaswani et al., 2017) encoders. The result from the DB derived from the output of DST is used for NLG inference. All of the expert layers share the same architecture. The input is ensembled by experts in the Soft-MoE layer for improving model capacity without the cost of efficiency. The model is fine-tuned by maximizing the likelihood of predicting the next token for NLU, DST, and NLG outputs.

084 dialogue models (Wu et al., 2020; Zhang et al., 2020; 117
 085 Peng et al., 2021; He et al., 2022b). However, their 118
 086 remarkable performance is at the cost of significant 119
 087 computational resources, especially as the sizes 120
 088 of PLMs continue to grow. Recently, parameter- 121
 089 efficient adapters raised that freeze the PLM while 122
 090 only allowing a small number of parameters up- 123
 091 dated for downstream models (Houlsby et al., 2019; 124
 092 Li and Liang, 2021; Lester et al., 2021), and have 125
 093 gained popularity in dialogue systems (Bang et al., 126
 094 2023; Wang et al., 2023). Nevertheless, the model 127
 095 capacity (i.e. number of parameters) is limited by 128
 096 the number of downstream models, and the addi- 129
 097 tion of adapters can become computationally expen- 130
 098 sive due to their sequential processing (Rücklé 131
 099 et al., 2020). We also argue that the issue of infer- 132
 100 ence time scaling with model complexity becomes 133
 101 more prominent considering the time sensitivity as- 134
 102 sociated with the deployment of dialogue systems. 135

103 To address these issues, we propose a Soft 136
 104 **Mixture-of-Expert Task-Oriented Dialogue** sys- 137
 105 tem (**SMETOD**) which scales the model capacity 138
 106 for diverse outputs of dialogue systems with signifi- 139
 107 cantly less training and inference cost. Specifically, 140
 108 we leverage Soft MoE (Puigcerver et al., 2023) to 141
 109 improve model capacity and leverage the effective- 142
 110 ness and performance of considerably larger mod- 143
 111 els with significantly lower computational costs. 144
 112 We present a task-oriented dialogue system as a 145
 113 multi-module end-to-end text generation to bridge 146
 114 the gap between traditional pipeline-based and end-
 115 to-end response generation systems, and optimize
 116 NLU, DST, and NLG, respectively, as in (Su et al.,

2022; Bang et al., 2021). We formulate NLU, DST, 117
 and NLG as the text generation problems, which 118
 take dialogue history sequence as model input and 119
 generate spans as the output. In the cases of NLG, 120
 we predict the DST output to obtain the database 121
 (DB) state, which becomes incorporated into its 122
 input. With T5-small (Raffel et al., 2020a) and T5- 123
 base (Raffel et al., 2020a) as the backbone PLM, 124
 we evaluate our method on MultiWOZ (Eric et al., 125
 2019; Zang et al., 2020) and NLU (Casanueva 126
 et al., 2020; Larson et al., 2019; Liu et al., 2019) 127
 datasets. We show that our method achieves signif- 128
 icant improvement in multi-domain DST on Multi- 129
 WOZ 2.1 and NLG on both benchmarks. 130

Our contribution is as follows: 131

- We propose SMETOD, a task-oriented dia- 132
 133
 134
 135
 136
 137
 138
 139
 140
 141
 142
 143
 144
 145
 146
 147
 148
 149
 150
 151
 152
 153
 154
 155
 156
 157
 158
 159
 160
 161
 162
 163
 164
 165
 166
 167
 168
 169
 170
 171
 172
 173
 174
 175
 176
 177
 178
 179
 180
 181
 182
 183
 184
 185
 186
 187
 188
 189
 190
 191
 192
 193
 194
 195
 196
 197
 198
 199
 200
 201
 202
 203
 204
 205
 206
 207
 208
 209
 210
 211
 212
 213
 214
 215
 216
 217
 218
 219
 220
 221
 222
 223
 224
 225
 226
 227
 228
 229
 230
 231
 232
 233
 234
 235
 236
 237
 238
 239
 240
 241
 242
 243
 244
 245
 246
 247
 248
 249
 250
 251
 252
 253
 254
 255
 256
 257
 258
 259
 260
 261
 262
 263
 264
 265
 266
 267
 268
 269
 270
 271
 272
 273
 274
 275
 276
 277
 278
 279
 280
 281
 282
 283
 284
 285
 286
 287
 288
 289
 290
 291
 292
 293
 294
 295
 296
 297
 298
 299
 300
 301
 302
 303
 304
 305
 306
 307
 308
 309
 310
 311
 312
 313
 314
 315
 316
 317
 318
 319
 320
 321
 322
 323
 324
 325
 326
 327
 328
 329
 330
 331
 332
 333
 334
 335
 336
 337
 338
 339
 340
 341
 342
 343
 344
 345
 346
 347
 348
 349
 350
 351
 352
 353
 354
 355
 356
 357
 358
 359
 360
 361
 362
 363
 364
 365
 366
 367
 368
 369
 370
 371
 372
 373
 374
 375
 376
 377
 378
 379
 380
 381
 382
 383
 384
 385
 386
 387
 388
 389
 390
 391
 392
 393
 394
 395
 396
 397
 398
 399
 400
 401
 402
 403
 404
 405
 406
 407
 408
 409
 410
 411
 412
 413
 414
 415
 416
 417
 418
 419
 420
 421
 422
 423
 424
 425
 426
 427
 428
 429
 430
 431
 432
 433
 434
 435
 436
 437
 438
 439
 440
 441
 442
 443
 444
 445
 446
 447
 448
 449
 450
 451
 452
 453
 454
 455
 456
 457
 458
 459
 460
 461
 462
 463
 464
 465
 466
 467
 468
 469
 470
 471
 472
 473
 474
 475
 476
 477
 478
 479
 480
 481
 482
 483
 484
 485
 486
 487
 488
 489
 490
 491
 492
 493
 494
 495
 496
 497
 498
 499
 500
 501
 502
 503
 504
 505
 506
 507
 508
 509
 510
 511
 512
 513
 514
 515
 516
 517
 518
 519
 520
 521
 522
 523
 524
 525
 526
 527
 528
 529
 530
 531
 532
 533
 534
 535
 536
 537
 538
 539
 540
 541
 542
 543
 544
 545
 546
 547
 548
 549
 550
 551
 552
 553
 554
 555
 556
 557
 558
 559
 560
 561
 562
 563
 564
 565
 566
 567
 568
 569
 570
 571
 572
 573
 574
 575
 576
 577
 578
 579
 580
 581
 582
 583
 584
 585
 586
 587
 588
 589
 590
 591
 592
 593
 594
 595
 596
 597
 598
 599
 600
 601
 602
 603
 604
 605
 606
 607
 608
 609
 610
 611
 612
 613
 614
 615
 616
 617
 618
 619
 620
 621
 622
 623
 624
 625
 626
 627
 628
 629
 630
 631
 632
 633
 634
 635
 636
 637
 638
 639
 640
 641
 642
 643
 644
 645
 646
 647
 648
 649
 650
 651
 652
 653
 654
 655
 656
 657
 658
 659
 660
 661
 662
 663
 664
 665
 666
 667
 668
 669
 670
 671
 672
 673
 674
 675
 676
 677
 678
 679
 680
 681
 682
 683
 684
 685
 686
 687
 688
 689
 690
 691
 692
 693
 694
 695
 696
 697
 698
 699
 700
 701
 702
 703
 704
 705
 706
 707
 708
 709
 710
 711
 712
 713
 714
 715
 716
 717
 718
 719
 720
 721
 722
 723
 724
 725
 726
 727
 728
 729
 730
 731
 732
 733
 734
 735
 736
 737
 738
 739
 740
 741
 742
 743
 744
 745
 746
 747
 748
 749
 750
 751
 752
 753
 754
 755
 756
 757
 758
 759
 760
 761
 762
 763
 764
 765
 766
 767
 768
 769
 770
 771
 772
 773
 774
 775
 776
 777
 778
 779
 780
 781
 782
 783
 784
 785
 786
 787
 788
 789
 790
 791
 792
 793
 794
 795
 796
 797
 798
 799
 800
 801
 802
 803
 804
 805
 806
 807
 808
 809
 810
 811
 812
 813
 814
 815
 816
 817
 818
 819
 820
 821
 822
 823
 824
 825
 826
 827
 828
 829
 830
 831
 832
 833
 834
 835
 836
 837
 838
 839
 840
 841
 842
 843
 844
 845
 846
 847
 848
 849
 850
 851
 852
 853
 854
 855
 856
 857
 858
 859
 860
 861
 862
 863
 864
 865
 866
 867
 868
 869
 870
 871
 872
 873
 874
 875
 876
 877
 878
 879
 880
 881
 882
 883
 884
 885
 886
 887
 888
 889
 890
 891
 892
 893
 894
 895
 896
 897
 898
 899
 900
 901
 902
 903
 904
 905
 906
 907
 908
 909
 910
 911
 912
 913
 914
 915
 916
 917
 918
 919
 920
 921
 922
 923
 924
 925
 926
 927
 928
 929
 930
 931
 932
 933
 934
 935
 936
 937
 938
 939
 940
 941
 942
 943
 944
 945
 946
 947
 948
 949
 950
 951
 952
 953
 954
 955
 956
 957
 958
 959
 960
 961
 962
 963
 964
 965
 966
 967
 968
 969
 970
 971
 972
 973
 974
 975
 976
 977
 978
 979
 980
 981
 982
 983
 984
 985
 986
 987
 988
 989
 990
 991
 992
 993
 994
 995
 996
 997
 998
 999
 1000

2 Preliminaries

Soft Mixture-of-Experts. Mixture-of-Experts (MoE)-based models have shown advantages in scaling model capacity without large increases in training or inference costs. There has been work on scaling sparsely activated MoE architectures. In the context of modern deep learning architectures, it was firstly found effective by Shazeer et al. (2017) by stacking MoE between LSTM (Hochreiter and Schmidhuber, 1997) and resulted in the state-of-the-art in language modeling and machine translation. Shazeer et al. introduced MoE Transformer where MoE layers are a substitute for the FFN layers (2018).

We adopt Soft-MoE (Puigcerver et al., 2023), which scales model capacity without the loss of fine-tuning efficiency and is fully differentiable and balanced compared to conventional efficient MoEs (Lepikhin et al., 2020; Fedus et al., 2022; Du et al., 2022; Zhou et al., 2022; Puigcerver et al., 2023). Specifically, it performs a soft assignment on experts to each input token, achieving similar training costs and much lower inference costs at a larger model capacity. We use $f(\cdot; \theta)$ to denote a mapping f associated with the parameter θ from the input sample to an output space. $\sigma(\cdot)$ is the Softmax function. We denote $\{f(\cdot; \theta_i)\}_{i=1}^m$ as m experts with identical architectures; their weights $\theta_1, \dots, \theta_m$ applied to individual tokens. Each expert has p slots, each of which is a weighted average of input. Slots in the same expert apply the same weights. Given input and output tokens $\mathbf{x} = \{x_1, \dots, x_l\}$ and $\mathbf{y} = \{y_1, \dots, y_l\}$ at the length l . Each expert will process p slots with parameters denoted as $\Psi = \{\psi^{(1)}, \dots, \psi^{(m \times p)}\}$. The input of experts, $\tilde{\mathbf{x}}$, is defined as the result of convex combinations of input tokens.

$$\tilde{\mathbf{x}}_j = (\sigma(\mathbf{x}\psi^{(j)}))^T \mathbf{x} \quad (1)$$

where j is the index of the slot in experts and $j \in [1, \dots, m \times p]$. The corresponding expert function is applied on each slot to obtain the output slots:

$$\tilde{\mathbf{y}}_j = f(\tilde{\mathbf{x}}_j; \theta_{\lfloor j/p \rfloor}) \quad (2)$$

Given $\tilde{\mathbf{y}} = \{\tilde{\mathbf{y}}_j\}_{j=1}^{m \times p}$, the output of Soft-MoE layer, y_i , is computed as a convex combination of all $(m \times p)$ output slots over the expert dimension (i.e. the rows of $\mathbf{x}\Psi$):

$$y_i = \sigma(x_i \Psi) \tilde{\mathbf{y}} \quad (3)$$

End-to-end task-oriented dialogue system. End-to-end learning was found effective in training and optimizing the map directly from input to output (Wen et al., 2017; Liu and Lane, 2018; Eric and Manning, 2017a; Williams et al., 2017). Later on, a lot of endeavor was given to fine-tuning pre-trained language models and adapting their generalization capacities for an end-to-end system of task-oriented dialogues (Budzianowski and Vulić, 2019; Casanueva et al., 2020; Mehri et al., 2020; Hosseini-Asl et al., 2020). In recent years, pre-trained task-oriented dialogue models have emerged as strong contenders, surpassing traditional fine-tuning approaches and showcasing competitive generalization capabilities, particularly in multi-objective scenarios (Wu et al., 2020; Zhang et al., 2020; Peng et al., 2021; He et al., 2022b). However, it’s worth noting that they require a large amount of dialogue data to train the backbone models and without an interface to optimize sub-modules.

Efficient transfer learning. To reduce the effort in tuning large PLMs and promote the scalability of model adaptation, there is a line of work that fixes the entire PLM and introduces a small number of new trainable parameters. Notable examples in this category include adapters (Houlsby et al., 2019; Pfeiffer et al., 2021; Karimi Mahabadi et al., 2021), prefix-tuning (Li and Liang, 2021) and prompt-tuning (Lester et al., 2021), etc. In-context learning prepends related task examples to condition on the generated dialogue states (Hu et al., 2022; Gupta et al., 2022; Venkateswaran et al., 2022). In end-to-end dialogue systems, a line of work prompts with specific text to generate desired outputs (Su et al., 2022) or injecting adapters to capture the knowledge of different functionalities (Bang et al., 2023; Mo et al., 2023). GPT-3 (Brown et al., 2020) and ChatGPT¹ are also successful and efficient open-domain dialogue systems. On the other hand, the MoE approach focuses on improving performance by efficiently scaling model sizes. Recent work on MoE develops more efficient routing implementations of Mixture-of-Experts in scaling language models (Lepikhin et al., 2020; Fedus et al., 2022; Du et al., 2022; Zhou et al., 2022; Puigcerver et al., 2023; Ma et al., 2018).

3 Method

We introduce **SMETOD**, a multi-objective dialogue system for NLU, DST, and NLG in task-

¹<https://chat.openai.com/chat>

oriented dialogues, scaling model capacities while maintaining computational efficiency with Soft MoE (Puigcerver et al., 2023). The overall architecture is illustrated in Figure 1.

3.1 Problem Formulation

We define the dialogue history $h = [u_1^{sys}, u_1^{usr}, \dots, u_t^{sys}, u_t^{usr}]$ as the concatenation of the system and user utterances in previous turns, where t is the number of current turns in the dialogue. h has all the dialogue history without the last system utterance, denoted as r . NLU outputs an I which is an intent or the API-name. The objective of DST is to output user goals, the tasks or purposes that the user wants to accomplish through the dialogue. user goals are typically represented as a set of pre-defined slot-value pairs that consist of the required information to query the dialogue system, i.e. $y_{API} = \{(s_1, v_1), \dots, (s_n, v_n)\}$, where n is the number of slot-value pairs. Finally, NLG will generate S with the previous output: $h + y_{DB} \rightarrow r$, where y_{DB} is the items in the database retrieved by y_{API} . Given a pair of training examples (x', y') , we elaborate x' and y' corresponding to different modules of the dialogue system in the following Table.

	x'	y'
NLU	h	I
DST	h	y_{API}
NLG	$h + y_{DB}$	r

3.2 Soft Mixture-of-Expert Layer

We implement the Soft-MoE layer to replace the second Feed-Forward Layer in each Transformer (Vaswani et al., 2017) Encoder block, as illustrated in Figure 1. Mathematically, we denote the output out the first Feed-Forward layer of the k -the encoder is $g(\cdot; \phi_k)$, then $\mathbf{x} = g(x'; \phi_k) \in \mathbb{R}^{l \times d_{ff}}$ in Eq. 1, denoting d_{ff} as the dimension between the first and second Feed-Forward layer and d as model’s hidden dimension, and l is the length of tokens. $\psi^{(j)} \in \mathbb{R}^{d_{ff}}$ is d_{ff} -dimensional vector of parameters corresponding to each slot of experts.

The mapping $f(\cdot; \theta_i)$ in Eq. 2 is simply a linear mapping corresponding to each expert, and p is the slots per expert having the same weights. Therefore, the output of the k -th encoder layer, $y'^{(k)}$, can

be represented as

$$y'^{(k)} = f(g(x'; \phi_k); \Theta_k, \Psi_k) \quad (4)$$

For fine-tuning, we replicate the pre-trained weights from the second Feed-Forward layer of encoders and assign them to each expert, leveraging the contextual learning abilities inherent in pre-trained models.

3.3 Training Objectives

We optimize the generation outputs of NLU, DST, NLG, respectively, following Su et al. (2022). Given a pair of training samples as (x', y') , the loss function is defined to maximize the log-likelihood of the token to predict given the current context:

$$\mathcal{L}_{\{NLU, DST, NLG\}} = -\frac{1}{l} \sum_{q=1}^l \log P(y'_q | y'_{<q}; x') \quad (5)$$

4 Experiment

4.1 Data

We evaluate our models for NLU on Banking77 (Casanueva et al., 2020), CLINC150 (Larson et al., 2019), and HWU64 (Liu et al., 2019); DST and NLG are evaluated on the task-oriented dialogue benchmarks MultiWOZ 2.1 (Eric et al., 2019) and MultiWOZ 2.2 (Zang et al., 2020). Banking77 contains 13,083 customer service queries labeled with 77 distinct intents for distinguishing between intents among queries related to similar tasks. CLINC150, consists of a comprehensive dataset comprising 23,700 examples, annotated with 150 intents across 10 distinct domains. HWU64 is collected from the home robot that has 25,716 examples for 64 intents spanning 21 domains.

MultiWOZ 2.1 (Eric et al., 2019) consists of multi-turn task-oriented dialogues across several domains, where 8,438 dialogues are for training and 1,0000 for dev and test. MultiWOZ 2.2 (Zang et al., 2020) improves MultiWOZ 2.1 by correcting annotation errors and adding dialogue act annotations. In MultiWOZ, the generation of response is not only related to the dialogue context but also grounded on the database (DB) state. The DB state is automatically retrieved from a pre-defined database using the generated dialogue states. SMETOD adopts a two-step approach during inference (Su et al., 2022; Bang et al., 2023). Firstly, it predicts the DST results to access the

Model	Banking77	HWU64	CLINC150
BERT-FIXED [◊] *	87.19	85.77	91.79
CONVBERT-DG	92.99	92.94	97.11
+Pre+Multi*			
CONVBERT	93.44	92.38	97.11
+Pre+Multi*			
BERT-TUNED [◊] *	93.66	92.10	96.93
CONVERT [◊] *	93.01	91.24	97.16
USE+CONVERT [◊] *	93.36	92.62	97.16
SPACE-2 [#] *	94.77	94.33*	97.80
SPACE-3*	94.94*	94.14	97.89
TOATOD _{small}	92.40	90.42	98.45
TOATOD _{base}	92.17	90.79	98.01
SMETOD _{small}	92.47	90.88	98.12
SMETOD _{base}	93.02	92.56	98.64

Table 1: Accuracy (%) on three intent prediction datasets with full-data experiments. [◊] comes from Casanueva et al. (2020). [#] are obtained from DialoGLUE leaderboard². All others are reported as in the original papers. Models with * are classification-based.

DB state. Subsequently, it utilizes the retrieved DB state and the current dialogue context to generate the NLG results.

4.2 Training & Inference Details

All models are fine-tuned respectively using PPTOD (Su et al., 2022), the pre-trained dialogue models based on T5-small (60M parameters) (Raffel et al., 2020b) and T5-base (220M parameters) (Raffel et al., 2020b), as the backbone. T5-small has 6 encoders and decoders with hidden size $d = 512$ and $d_{ff} = 2048$. While T5-large has 12 encoders and decoders and $d = 768$, $d_{ff} = 3072$. For models’ architecture, we replace the second Feed-Forward layer in all encoder blocks with the illustrated Soft-MoE layers, and copy pre-trained weights to each expert in the Soft-MoE layers. We augment T5 with 8 experts and 2 slots per expert for DST, and 16 experts with 2 slots per expert for NLU and NLG.

We fine-tuned all model parameters on the full-shot training setting. The linear combination weights in Soft-MoE layer are initialized by Kaiming initialization (He et al., 2015). The initial learning rate is set to 0.001 for NLU, and 0.0001 for DST, NLG, respectively. We use the Adafactor (Shazeer and Stern, 2018) optimizer and the training batch size is set to 64 on Nvidia A10 GPUs. We tried a wide range of learning rates from $1e-2$ to $1e-6$ then set the initial training rate to $1e-4$ in all training. Our code is developed based on *Soft-Mixture-of-*

*Experts*³ and *TOATOD*⁴. Code repository will be released to the public soon.

Because different batch sizes will result in different padded lengths, inference results are slightly changed by batch sizes due to Softmax over input tokens in the Soft-MoE layer. We make inferences on several selected batch sizes and report average scores. We found out that different batch sizes in our experiments have negligible influence on the inference results⁵.

5 Results & Discussion

We show the effectiveness of our models on NLU (Sec. 5.1), DST (Sec. 5.2), and NLG (Sec. 5.3) in task-oriented dialogue systems compared to plenty of strong baselines. In the experiments, we fine-tune SMETOD using the small and base versions of PPTOD (Su et al., 2022), which continues pre-training T5 (Raffel et al., 2020b) on large dialogue corpora, as the start point. We observe that SMETOD is state-of-the-art on NLU and DST and comparable with existing baselines on NLG. We also study the improvement of efficiency with SMETOD (Sec. 5.4). In Sec. 5.5, we investigate model performance when the Soft-MoE layers are in different architectures.

5.1 Intent Prediction

The goal of intent prediction, known as NLU in a task-oriented dialogue system, is to identify the user’s intention based on the user’s utterance. We conduct experiments on three benchmarks: Banking77 (Casanueva et al., 2020), CLINC150 (Larson et al., 2019), and HWU64 (Liu et al., 2019). We report Accuracy (%) of predicting an intention correctly for evaluation.

5.1.1 Baselines

Baselines have a wide range from BERT-based models: CONVBERT (Mehri et al., 2020), CONVERT (Casanueva et al., 2020), UniLM-based models: SPACE-2 (He et al., 2022a), SPACE-3 (He et al., 2022b), to T5-based TOATOD (Bang et al., 2023). All baseline models utilizing BERT and UniLM follow a classification-based approach, employing a classifier featuring a Softmax layer to make predictions from a predefined set of intents.

³<https://github.com/fkodom/soft-mixture-of-experts.git>

⁴<https://github.com/sogang-isds/TOATOD.git>

⁵We conducted a hypothesis test and found out p -value < 0.01 for scores changed by batch size. Statistics are summarized in Appendix A, Table 5.

Model	Pre-Trained Model	MultiWOZ2.1	MultiWOZ2.2
TRADE	-	45.6	45.4
TripPy	BERT-base	55.29	-
TripPy+SaCLog	BERT-base	60.61	-
CONVBERT-DG	BERT-base	55.29	-
SimpleTOD	DistilGPT-2	55.76	-
SOLOIST	GPT-2	56.85	-
AG-DST	PLATO-2	57.26	57.26
UniLM [‡]	UniLM	54.25	54.25
SPACE-3	UniLM	57.50	57.50
PPTOD _{base}	T5-base	57.10	-
PPTOD _{large}	T5-large	57.45	-
D3ST _{base}	T5-base	54.2	56.1
D3ST _{large}	T5-large	54.5	54.2
D3ST _{XXL}	T5-XXL	57.80	58.7
T5DST _{+desc}	T5-base	56.66	57.6
TOATOD _{small} [†]	T5-small	59.49	59.33
TOATOD _{base} [†]	T5-base	59.51	60.02
SMETOD _{small}	T5-small	59.69	59.60
SMETOD _{base}	T5-base	60.36	60.08

Table 2: Joint Goal Accuracy (%) for DST on MultiWOZ 2.1 and 2.2. Results with [‡] are from He et al.(2022b). [†] represents the results of our re-implementation. All others are reported as in the original papers.

5.1.2 Evaluation Results

Table 1 shows that our approaches perform state-of-art on CLINC150, which has the most number of intent types. On the other two benchmarks, our approaches have the highest accuracy compared to other generation-based approaches. Classification-based approaches are better which may benefit from smaller numbers of intents to choose from. Compared to classification models, SMETOD copes with the classification task as a generation problem by directly generating the text label. Therefore, when adapting to a new classification task, SMETOD is more scalable to new domains and tasks and can predict intents that are not in the ontology.

5.2 Dialogue State Tracking

As a crucial component in task-oriented dialogue systems, DST determines the user goals based on the history of dialogue turns. For the evaluation of DST models, we use joint goal accuracy (JGA) which is the average accuracy of predicting all slot-values for the current turn correctly.

5.2.1 Baselines

In Table 2, we compare SMETOD with a wide range of classification-based approaches: TRADE (Wu et al., 2019), TripPy (Heck et al., 2020), TripPy + SaCLog (Dai et al., 2021), CONVBERT-DG (Mehri et al., 2020), Simple-

TOD (Hosseini-Asl et al., 2020), SOLOIST (Peng et al., 2021), AG-DST (Tian et al., 2021), SPACE-3 (He et al., 2022b), and generation-based approaches: PPTOD (Su et al., 2022), D3ST (Zhao et al., 2022), T5DST (Lee et al., 2021), and TOATOD (Bang et al., 2023).

5.2.2 Evaluation Results

Compared to other approaches, SMETOD obtains state-of-the-art JGA on MultiWOZ 2.1 and 2.2 among all generation-based approaches. Our model is more flexible to generate slot-value pairs while classification-based models are limited to the pre-defined ontology. The results show that our model can benefit from not only the transfer learning capacities of pre-trained models but also the improvement of model size.

5.3 End-to-End Response Generation

End-to-end dialogue response generation, aiming at evaluating the model in the most realistic, fully end-to-end setting, where the generated dialogue states are used for the database search and response generation (Hosseini-Asl et al., 2020; Su et al., 2022), is NLG in task-oriented dialogue system. Our models evaluated on MultiWOZ generates responses not only related to the dialogue history but also grounded on the database (DB) state.

Model	Backbone	MultiWOZ2.1				MultiWOZ2.2			
		Inform	Success	BLEU	Combined	Inform	Success	BLEU	Combined
DOTS	BERT-base	86.65	74.18	15.90	96.32	-	-	-	-
DiactTOD	S-BERT	-	-	-	-	89.5	84.2	17.5	104.4
SimpleTOD	DistilGPT-2	85.00	70.50	15.23	92.98	-	-	-	-
SOLOIST	GPT-2	-	-	-	-	82.3	72.4	13.6	90.9
UBAR ^Δ	GPT-2	95.70	81.80	16.50	105.25	83.4	70.3	17.6	94.4
MinTL ^Δ	BART _{large}	-	-	-	-	73.7	65.4	19.4	89.0
RewardNet ^Δ	BART _{large}	-	-	-	-	87.6	81.5	17.6	102.2
GALAXY	UniLM	95.30	86.20	20.01	110.76	85.4	75.7	19.64	100.2
PPTOD _{base}	T5-base	87.09	79.08	19.17	102.26	-	-	-	-
MTTOD [‡]	T5-base	90.99	82.08	19.68	106.22	85.9	76.5	19.0	100.2
RSTOD [‡]	T5-small	93.50	84.70	19.24	108.34	83.5	75.0	18.0	97.3
TOATOD _{small}	T5-small	92.10	80.40	18.29	104.54	85.80	74.00	18.00	97.90
TOATOD _{base}	T5-base	97.00	87.40	17.12	109.32	90.00	79.80	17.04	101.94
KRLS	T5-base	-	-	-	-	89.2	80.3	19.0	103.8
SMETOD _{small}	T5-small	92.50	74.00	16.89	100.14	89.6	76.2	17.1	100.1
SMETOD _{base}	T5-base	92.30	78.80	16.88	102.43	89.0	76.0	17.6	99.7

Table 3: Evaluation of NLG on Inform, Success, BLEU, and Combined Scores, where Combined = (Inform + Success) × 0.5 + BLEU. [‡] means the NLG results on MultiWOZ 2.1 is from [Cholakov and Kolev \(2022\)](#). All other results are from MultiWOZ leaderboards⁶. ^Δ shows models that require oracle dialogue states for prediction.

5.3.1 Metrics

For evaluation, we follow the individual and combined metrics in [Hosseini-Asl et al. \(2020\)](#): Inform, Success, and BLEU, and Combined score which is defined as Combined = (Inform + Success) × 0.5 + BLEU. Specifically, Inform rate measures the correctness of entities in the response. Success rate success rate assesses attribute fulfillment requested by user. BLEU score is used to measure the fluency of the generated responses.

5.3.2 Baselines

In Table 3, we compare our model with several strong baselines: DOTS ([Jeon and Lee, 2021](#)), DiactTOD ([Wu et al., 2023](#)), SimpleTOD ([Hosseini-Asl et al., 2020](#)), SOLOIST ([Peng et al., 2021](#)), UBAR ([Yang et al., 2021](#)), MinTL ([Lin et al., 2020](#)), RewardNet ([Feng et al., 2023](#)), GALAXY ([He et al., 2022c](#)), PPTOD ([Su et al., 2022](#)), RSTOD ([Cholakov and Kolev, 2022](#)), MTTOD ([Lee, 2021](#)), TOATOD ([Bang et al., 2023](#)), KRLS ([Xiao Yu, 2022](#)).

5.3.3 Evaluation Results

On both MultiWOZ 2.1 and 2.2 datasets, SMETOD performs, though not the best, comparable to T5-based models except TOATOD_{base}. We hypothesize that metrics hinder each other from being improved together and may require a mech-

anism to promote performance towards specific metrics, for example, REINFORCE ([Sutton et al., 1999](#)). Besides, we observe that only replacing the Feed-Forward layer in Transformer encoders as in [Puigcerver et al. \(2023\)](#) without copying weights to experts doesn't generate the best results in our dialogue system. It might be because their implementation requires a large amount of data to pre-train, which is inappropriate in the task-oriented scenario. It demonstrates that by duplicating pre-trained weights and fine-tuning, SMETOD optimizes well for DST and NLG, respectively, maintaining the prior knowledge learned from the pre-trained model.

Model	Small↓	Base↓
PPTOD	1×	3.163×
TOATOD	1.116×	3.519×
SMETOD	1.005×	3.095×

Table 4: Comparison of the inference time with small and base-size models of PPTOD and TOATOD for NLG on MultiWOZ 2.1. All models are experimented with 5 same and randomly sampled batch sizes. Average time is reported. ↓: Smaller is better.

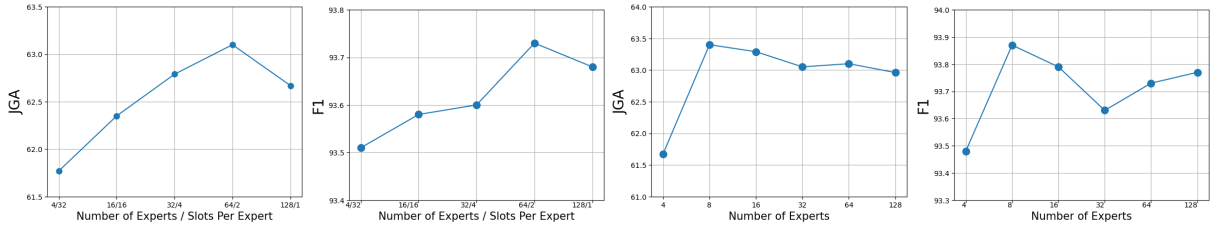


Figure 2: **(Left)** Performance of SMETOD as a function of the number of experts, for models with a fixed number of experts \times slots-per-expert. **(Right)** Performance of SMETOD trained with increased experts and 2 slots per expert. JGA and F1 scores are on MultiWOZ 2.1 dev set for DST.

5.4 Time Complexity Analysis

According to Puigcerver et al. (2023), the time complexity of the Soft-MoE layer can be reduced to $O(l^2d + lk)$, given input token length l , model hidden dimension d , and the cost of applying an expert per token $O(k)$. Thus, the time complexity is constant and the same as the single-headed self-attention cost by increasing the number of experts m and scaling slots per expert $p = O(l/m)$ accordingly, which will not become a bottleneck in Transformer.

We show in Table 4 that SMETOD could make inferences without bringing about much extra time. SMETOD_{small} is 3.5 times larger than PPTOD and TOATOD while achieving a similar inference speed as the former. Our SMETOD_{base} has even less inference time while its model size is 4 times of PPTOD_{base}. It proves that we can achieve much better scaling while cost is roughly constant (Puigcerver et al., 2023), with the benefit of improved performance.

5.5 Impact of Expert Numbers

We investigate the impact of expert and slot numbers in our models on the development set of MultiWOZ 2.1 for DST as illustrated in Figure 2. First, we fix the total number of slots to 128 and vary expert numbers {4, 16, 32, 64, 128} by scaling slot numbers per expert. Results suggest the best configuration is 64 experts and 2 slots per expert. Then, we set the number of slots per expert to one and evaluate performance with regard to the number of experts. The number of experts 8 and 16 perform better than others. It should be mentioned that the model size scales with increasing expert numbers only. Meanwhile, we observe performance is not always increasing with the number of experts, indicating there is a trade-off between model size and the amount of training data.

6 Conclusion

We propose an efficient fine-tuning approach based on Soft-MoE to satisfy requirements on diverse outputs in task-oriented dialogue systems. We demonstrate that incorporating Soft-MoE to our dialogue system achieves remarkable success on MultiWOZ baselines and optimizes outputs of each submodule, showing it powerful technique for task-oriented dialogue systems with better scaling performance while maintaining time efficiency.

7 Limitations

Limitations related to adopting Soft-MoE: This work is a practice of leveraging Soft-MoE (Puigcerver et al., 2023) in downstream models with supervised, while the original practice requires unsupervised pre-training. We consider per-taining experts on larger dialogue corpus, for example, Lin et al. (2021); Hu et al. (2022) for better generality performance in the future. Furthermore, we didn't evaluate our approach to NLP datasets which have more diverse example lengths. Unlike Soft-MoE used in computer vision, the weights over tokens are inconstant due to the variety of length of input tokens, which leads to inconsistent inference with different batch sizes. Although we observe negligible influence in our experiments, variations of lengths require further study. We should also have experimented with more expert numbers and investigated the performance on NLG as well to study how performance is improved with model size. Last, scaling up model sizes requires a lot of computational memory.

Limitations related to datasets: DST and NLG evaluations are on MultiWOZ, which are English and have limited domains. More generalized and larger-scale dialogue corpus need considering, such as DialoGLUE (Mehri et al., 2020), SGD (Lee

et al., 2022), or multi-lingual datasets (Ding et al., 2021). NLU evaluations are only on single-utterance benchmarks, CamRest676 (Quan and Xiong, 2019), In-Car Assistant (Eric and Manning, 2017b).

Limitations related to training time: Recently, Adapters and prompt approaches have been proposed that update fewer parameters in models compared with our fine-tuning approaches. Although we didn’t observe longer training time explicitly compared to adapter-based models with similar sizes, empirical study on this issue is not covered in this work. We have shown in Sec. 5.4 that the forward pass of our approach is faster. It has been shown that original adapters should backpropagate through the entire model only except the first components (Rücklé et al., 2020). Moreover, we argue that performance and inference efficiency are more important regarding the deployment of task-oriented dialogue systems.

Limitations related to GPT3 or ChatGPT (LLM) as baselines: We did not include evaluation with the above models due to the following reasons. First, we consider the generation problem in this paper to generate diverse outputs given the same input. The quality of prompts will have a significant impact on LLM results, making it hard to make a fair comparison. Second, our training is in full-shot scenarios, while GPT3 or ChatGPT is usually considered as a zero-shot or few-shot baseline. Last, there is a high probability that LLMs have contaminated public benchmarks used in this paper.

8 Potential Risks

Using public dialogue benchmarks introduces the potential for biases stemming from the data collection method. Models trained on such datasets might encounter challenges when attempting to generalize to real-world scenarios or specific domains, as the data may not accurately represent these situations. Additionally, public dialogue datasets frequently lack essential context or metadata, rendering it difficult to comprehend the circumstances surrounding the conversations.

In our approach, we also rely on open-source code repositories. However, these repositories can present issues related to security vulnerabilities and compatibility. Furthermore, their often incomplete documentation can pose additional hurdles for further development. Given the absence of reliable support or comprehensive documentation, these

factors can impede troubleshooting and hinder the overall development process.

References

- Namo Bang, Jeehyun Lee, and Myoung-Wan Koo. 2023. Task-optimized adapters for an end-to-end task-oriented dialogue system. *arXiv preprint arXiv:2305.02468*.
- Yejin Bang, Nayeon Lee, Etsuko Ishii, Andrea Madotto, and Pascale Fung. 2021. [Assessing political prudence of open-domain chatbots](#). In *Proceedings of the 22nd Annual Meeting of the Special Interest Group on Discourse and Dialogue*, pages 548–555, Singapore and Online. Association for Computational Linguistics.
- Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. 2020. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901.
- Paweł Budzianowski and Ivan Vulić. 2019. Hello, it’s gpt-2—how can i help you? towards the use of pre-trained language models for task-oriented dialogue systems. *arXiv preprint arXiv:1907.05774*.
- Iñigo Casanueva, Tadas Temčinas, Daniela Gerz, Matthew Henderson, and Ivan Vulić. 2020. Efficient intent detection with dual sentence encoders. *arXiv preprint arXiv:2003.04807*.
- Radostin Cholakov and Todor Kolev. 2022. Efficient task-oriented dialogue systems with response selection as an auxiliary task. *arXiv preprint arXiv:2208.07097*.
- Yinpei Dai, Hangyu Li, Yongbin Li, Jian Sun, Fei Huang, Luo Si, and Xiaodan Zhu. 2021. Preview, attend and review: Schema-aware curriculum learning for multi-domain dialog state tracking. *arXiv preprint arXiv:2106.00291*.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2018. [Bert: Pre-training of deep bidirectional transformers for language understanding](#). *arXiv preprint arXiv:1810.04805*.
- Bosheng Ding, Junjie Hu, Lidong Bing, Sharifah Mahani Aljunied, Shafiq Joty, Luo Si, and Chunyan Miao. 2021. Globalwoz: Globalizing multiwoz to develop multilingual task-oriented dialogue systems. *arXiv preprint arXiv:2110.07679*.
- Li Dong, Nan Yang, Wenhui Wang, Furu Wei, Xiaodong Liu, Yu Wang, Jianfeng Gao, Ming Zhou, and Hsiao-Wuen Hon. 2019. [Unified language model pre-training for natural language understanding and generation](#). *Advances in neural information processing systems*, 32.

677	Nan Du, Yanping Huang, Andrew M Dai, Simon Tong,	and generation. In <i>Proceedings of the 45th Inter-</i>	733
678	Dmitry Lepikhin, Yuanzhong Xu, Maxim Krikun,	<i>national ACM SIGIR Conference on Research and</i>	734
679	Yanqi Zhou, Adams Wei Yu, Orhan Firat, et al. 2022.	<i>Development in Information Retrieval</i> , pages 187–	735
680	Glam: Efficient scaling of language models with	200.	736
681	mixture-of-experts. In <i>International Conference on</i>		
682	<i>Machine Learning</i> , pages 5547–5569. PMLR.		
683	Mihail Eric, Rahul Goel, Shachi Paul, Abhishek Sethi,	Wanwei He, Yinpei Dai, Yinhe Zheng, Yuchuan Wu,	737
684	Sanchit Agarwal, Shuyang Gao, and Dilek Hakkani-	Zheng Cao, Dermot Liu, Peng Jiang, Min Yang, Fei	738
685	Tur. 2019. Multiwoz 2.1: Multi-domain dialogue	Huang, Luo Si, et al. 2022c. Galaxy: A generative	739
686	state corrections and state tracking baselines. <i>arXiv</i>	pre-trained model for task-oriented dialog with semi-	740
687	<i>preprint arXiv:1907.01669</i> .	supervised learning and explicit policy injection. In	741
		<i>Proceedings of the AAAI conference on artificial in-</i>	742
		<i>telligence</i> , volume 36, pages 10749–10757.	743
688	Mihail Eric and Christopher D Manning. 2017a. A	Michael Heck, Carel van Niekerk, Nurul Lubis, Chris-	744
689	copy-augmented sequence-to-sequence architecture	tian Geishauser, Hsien-Chin Lin, Marco Moresi, and	745
690	gives good performance on task-oriented dialogue.	Milica Gašić. 2020. Trippy: A triple copy strategy	746
691	<i>arXiv preprint arXiv:1701.04024</i> .	for value independent neural dialog state tracking.	747
		<i>arXiv preprint arXiv:2005.02877</i> .	748
692	Mihail Eric and Christopher D Manning. 2017b. Key-	Sepp Hochreiter and Jürgen Schmidhuber. 1997. Long	749
693	value retrieval networks for task-oriented dialogue.	short-term memory. <i>Neural computation</i> , 9(8):1735–	750
694	<i>arXiv preprint arXiv:1705.05414</i> .	1780.	751
695	William Fedus, Barret Zoph, and Noam Shazeer. 2022.	Ehsan Hosseini-Asl, Bryan McCann, Chien-Sheng Wu,	752
696	Switch transformers: Scaling to trillion paramete-	Semih Yavuz, and Richard Socher. 2020. A simple	753
697	ter models with simple and efficient sparsity. <i>The</i>	language model for task-oriented dialogue. <i>Advances</i>	754
698	<i>Journal of Machine Learning Research</i> , 23(1):5232–	<i>in Neural Information Processing Systems</i> , 33:20179–	755
699	5270.	20191.	756
700	Yihao Feng, Shentao Yang, Shujian Zhang, Jianguo	Neil Houlsby, Andrei Giurgiu, Stanislaw Jastrzebski,	757
701	Zhang, Caiming Xiong, Mingyuan Zhou, and Huan	Bruna Morrone, Quentin De Laroussilhe, Andrea	758
702	Wang. 2023. Fantastic rewards and how to tame them:	Gesmundo, Mona Attariyan, and Sylvain Gelly. 2019.	759
703	A case study on reward learning for task-oriented	Parameter-efficient transfer learning for nlp . In <i>In-</i>	760
704	dialogue systems. <i>arXiv preprint arXiv:2302.10342</i> .	<i>ternational Conference on Machine Learning</i> , pages	761
705	D. Goddeau, H. Meng, J. Polifroni, S. Seneff, and	2790–2799. PMLR.	762
706	S. Busayapongchai. 1996. A form-based dialogue	Yushi Hu, Chia-Hsuan Lee, Tianbao Xie, Tao Yu,	763
707	manager for spoken language applications . In <i>Pro-</i>	Noah A Smith, and Mari Ostendorf. 2022. In-context	764
708	<i>ceeding of Fourth International Conference on Spo-</i>	learning for few-shot dialogue state tracking. <i>arXiv</i>	765
709	<i>ken Language Processing. ICSLP '96</i> , volume 2,	<i>preprint arXiv:2203.08568</i> .	766
710	pages 701–704 vol.2.		
711	Raghav Gupta, Harrison Lee, Jeffrey Zhao, Abhinav	Hyunmin Jeon and Gary Geunbae Lee. 2021. Domain	767
712	Rastogi, Yuan Cao, and Yonghui Wu. 2022. Show,	state tracking for a simplified dialogue system. <i>arXiv</i>	768
713	don't tell: Demonstrations outperform descriptions	<i>preprint arXiv:2103.06648</i> .	769
714	for schema-guided task-oriented dialogue. <i>arXiv</i>		
715	<i>preprint arXiv:2204.04327</i> .	Rabeeh Karimi Mahabadi, James Henderson, and Se-	770
		bastian Ruder. 2021. Compacter: Efficient low-rank	771
716	Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian	hypercomplex adapter layers. <i>Advances in Neural</i>	772
717	Sun. 2015. Delving deep into rectifiers: Surpassing	<i>Information Processing Systems</i> , 34:1022–1035.	773
718	human-level performance on imagenet classification.	Stefan Larson, Anish Mahendran, Joseph J. Peper,	774
719	In <i>Proceedings of the IEEE International Conference</i>	Christopher Clarke, Andrew Lee, Parker Hill,	775
720	<i>on Computer Vision (ICCV)</i> .	Jonathan K. Kummerfeld, Kevin Leach, Michael A.	776
721	Wanwei He, Yinpei Dai, Binyuan Hui, Min Yang, Zheng	Laurenzano, Lingjia Tang, and Jason Mars. 2019. An	777
722	Cao, Jianbo Dong, Fei Huang, Luo Si, and Yong-	evaluation dataset for intent classification and out-of-	778
723	bin Li. 2022a. SPACE-2: Tree-structured semi-	scope prediction . In <i>Proceedings of the 2019 Confer-</i>	779
724	supervised contrastive pre-training for task-oriented	<i>ence on Empirical Methods in Natural Language Pro-</i>	780
725	dialog understanding . In <i>Proceedings of the 29th</i>	<i>cessing and the 9th International Joint Conference</i>	781
726	<i>International Conference on Computational Linguis-</i>	<i>on Natural Language Processing (EMNLP-IJCNLP)</i> ,	782
727	<i>tics</i> , pages 553–569, Gyeongju, Republic of Korea.	pages 1311–1316, Hong Kong, China. Association	783
728	International Committee on Computational Linguis-	for Computational Linguistics.	784
729	tics.		
730	Wanwei He, Yinpei Dai, Min Yang, Jian Sun, Fei Huang,	Chia-Hsuan Lee, Hao Cheng, and Mari Ostendorf. 2021.	785
731	Luo Si, and Yongbin Li. 2022b. Unified dialog model	Dialogue state tracking with a language model using	786
732	pre-training for task-oriented dialog understanding	schema-driven prompting . In <i>Proceedings of the</i>	787
		<i>2021 Conference on Empirical Methods in Natural</i>	788

789		<i>Language Processing</i> , pages 4937–4949, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.	
790			
791			
792	Harrison Lee, Raghav Gupta, Abhinav Rastogi, Yuan Cao, Bin Zhang, and Yonghui Wu. 2022. Sgd-x: A benchmark for robust generalization in schema-guided dialogue systems. In <i>Proceedings of the AAAI Conference on Artificial Intelligence</i> , volume 36, pages 10938–10946.		
793			
794			
795			
796			
797			
798	Hwaran Lee, Jinsik Lee, and Tae-Yoon Kim. 2019. Sumbt: Slot-utterance matching for universal and scalable belief tracking. <i>arXiv preprint arXiv:1907.07421</i> .		
799			
800			
801			
802	Yohan Lee. 2021. Improving end-to-end task-oriented dialog system with a simple auxiliary task. In <i>Findings of the Association for Computational Linguistics: EMNLP 2021</i> , pages 1296–1303.		
803			
804			
805			
806	Dmitry Lepikhin, HyoukJoong Lee, Yuanzhong Xu, Dehao Chen, Orhan Firat, Yanping Huang, Maxim Krikun, Noam Shazeer, and Zhifeng Chen. 2020. Gshard: Scaling giant models with conditional computation and automatic sharding. <i>arXiv preprint arXiv:2006.16668</i> .		
807			
808			
809			
810			
811			
812	Brian Lester, Rami Al-Rfou, and Noah Constant. 2021. The power of scale for parameter-efficient prompt tuning . <i>arXiv preprint arXiv:2104.08691</i> .		
813			
814			
815	Xiang Lisa Li and Percy Liang. 2021. Prefix-tuning: Optimizing continuous prompts for generation. <i>arXiv preprint arXiv:2101.00190</i> .		
816			
817			
818	Zhaojiang Lin, Bing Liu, Andrea Madotto, Seungwhan Moon, Paul Crook, Zhenpeng Zhou, Zhiguang Wang, Zhou Yu, Eunjoon Cho, Rajen Subba, et al. 2021. Zero-shot dialogue state tracking via cross-task transfer. <i>arXiv preprint arXiv:2109.04655</i> .		
819			
820			
821			
822			
823	Zhaojiang Lin, Andrea Madotto, Genta Indra Winata, and Pascale Fung. 2020. Mintl: Minimalist transfer learning for task-oriented dialogue systems. <i>arXiv preprint arXiv:2009.12005</i> .		
824			
825			
826			
827	Bing Liu and Ian Lane. 2018. End-to-end learning of task-oriented dialogs. In <i>Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Student Research Workshop</i> , pages 67–73.		
828			
829			
830			
831			
832	Xingkun Liu, Arash Eshghi, Pawel Swietojanski, and Verena Rieser. 2019. Benchmarking natural language understanding services for building conversational agents . <i>arXiv preprint arXiv:1903.05566</i> .		
833			
834			
835			
836	Jiaqi Ma, Zhe Zhao, Xinyang Yi, Jilin Chen, Lichan Hong, and Ed H Chi. 2018. Modeling task relationships in multi-task learning with multi-gate mixture-of-experts. In <i>Proceedings of the 24th ACM SIGKDD international conference on knowledge discovery & data mining</i> , pages 1930–1939.		
837			
838			
839			
840			
841			
	Shikib Mehri, Mihail Eric, and Dilek Hakkani-Tur. 2020. Dialogue: A natural language understanding benchmark for task-oriented dialogue. <i>arXiv preprint arXiv:2009.13570</i> .		842 843 844 845
	Yunho Mo, Joon Yoo, and Sangwoo Kang. 2023. Parameter-efficient fine-tuning method for task-oriented dialogue systems. <i>Mathematics</i> , 11(14):3048.		846 847 848 849
	Nikola Mrkšić, Diarmuid O Séaghdha, Tsung-Hsien Wen, Blaise Thomson, and Steve Young. 2016. Neural belief tracker: Data-driven dialogue state tracking. <i>arXiv preprint arXiv:1606.03777</i> .		850 851 852 853
	Baolin Peng, Chunyuan Li, Jinchao Li, Shahin Shayan-deh, Lars Liden, and Jianfeng Gao. 2021. Soloist: Building task bots at scale with transfer learning and machine teaching. <i>Transactions of the Association for Computational Linguistics</i> , 9:807–824.		854 855 856 857 858
	Jonas Pfeiffer, Aishwarya Kamath, Andreas Rücklé, Kyunghyun Cho, and Iryna Gurevych. 2021. AdapterFusion: Non-destructive task composition for transfer learning . In <i>Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume</i> , pages 487–503, Online. Association for Computational Linguistics.		859 860 861 862 863 864 865 866
	Joan Puigcerver, Carlos Riquelme, Basil Mustafa, and Neil Houlsby. 2023. From sparse to soft mixtures of experts. <i>arXiv preprint arXiv:2308.00951</i> .		867 868 869
	Jun Quan and Deyi Xiong. 2019. Effective data augmentation approaches to end-to-end task-oriented dialogue. In <i>2019 International Conference on Asian Language Processing (IALP)</i> , pages 47–52. IEEE.		870 871 872 873
	Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. 2019. Language models are unsupervised multitask learners. <i>OpenAI blog</i> , 1(8):9.		874 875 876 877
	Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J. Liu. 2020a. Exploring the limits of transfer learning with a unified text-to-text transformer . <i>Journal of Machine Learning Research</i> , 21(140):1–67.		878 879 880 881 882 883
	Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, Peter J Liu, et al. 2020b. Exploring the limits of transfer learning with a unified text-to-text transformer. <i>J. Mach. Learn. Res.</i> , 21(140):1–67.		884 885 886 887 888
	Andreas Rücklé, Gregor Geigle, Max Glockner, Tilman Beck, Jonas Pfeiffer, Nils Reimers, and Iryna Gurevych. 2020. Adapterdrop: On the efficiency of adapters in transformers. <i>arXiv preprint arXiv:2010.11918</i> .		889 890 891 892 893
	Noam Shazeer, Youlong Cheng, Niki Parmar, Dustin Tran, Ashish Vaswani, Penporn Koanantakool, Peter Hawkins, HyoukJoong Lee, Mingsheng Hong, Cliff		894 895 896

897	Young, et al. 2018. Mesh-tensorflow: Deep learning for supercomputers. <i>Advances in neural information processing systems</i> , 31.	951
898		952
899		953
900	Noam Shazeer, *Azalia Mirhoseini, *Krzysztof	954
901	Maziarz, Andy Davis, Quoc Le, Geoffrey Hinton,	955
902	and Jeff Dean. 2017. Outrageously large neural networks: The sparsely-gated mixture-of-experts layer .	956
903	In <i>International Conference on Learning Representations</i> .	
904		
905		
906	Noam Shazeer and Mitchell Stern. 2018. Adafactor:	
907	Adaptive learning rates with sublinear memory cost.	
908	In <i>International Conference on Machine Learning</i> ,	
909	pages 4596–4604. PMLR.	
910	Pei-Hao Su, Milica Gasic, Nikola Mrksic, Lina Rojas-	
911	Barahona, Stefan Ultes, David Vandyke, Tsung-	
912	Hsien Wen, and Steve Young. 2016. Continuously	
913	learning neural dialogue management. <i>arXiv preprint</i>	
914	<i>arXiv:1606.02689</i> .	
915	Yixuan Su, Lei Shu, Elman Mansimov, Arshit Gupta,	
916	Deng Cai, Yi-An Lai, and Yi Zhang. 2022. Multi-task	
917	pre-training for plug-and-play task-oriented dialogue	
918	system .	
919	Richard S Sutton, David McAllester, Satinder Singh,	
920	and Yishay Mansour. 1999. Policy gradient methods	
921	for reinforcement learning with function approxima-	
922	tion. <i>Advances in neural information processing</i>	
923	<i>systems</i> , 12.	
924	Ryuichi Takanobu, Qi Zhu, Jinchao Li, Baolin Peng,	
925	Jianfeng Gao, and Minlie Huang. 2020. Is your goal-	
926	oriented dialog model performing really well? empiri-	
927	cal analysis of system-wise evaluation. <i>arXiv</i>	
928	<i>preprint arXiv:2005.07362</i> .	
929	Blaise Thomson and Steve Young. 2010. Bayesian	
930	update of dialogue state: A pomdp framework for	
931	spoken dialogue systems. <i>Computer Speech & Lan-</i>	
932	<i>guage</i> , 24(4):562–588.	
933	Xin Tian, Liankai Huang, Yingzhan Lin, Siqi Bao,	
934	Huang He, Yunyi Yang, Hua Wu, Fan Wang, and	
935	Shuqi Sun. 2021. Amendable generation for dialogue	
936	state tracking. <i>arXiv preprint arXiv:2110.15659</i> .	
937	Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob	
938	Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz	
939	Kaiser, and Illia Polosukhin. 2017. Attention is all	
940	you need. <i>Advances in neural information processing</i>	
941	<i>systems</i> , 30.	
942	Praveen Venkateswaran, Evelyn Duesterwald, and	
943	Vatche Isahagian. 2022. District: Dialogue state	
944	tracking with retriever driven in-context tuning .	
945	<i>arXiv preprint arXiv:2212.02851</i> .	
946	Qingyue Wang, Liang Ding, Yanan Cao, Yibing Zhan,	
947	Zheng Lin, Shi Wang, Dacheng Tao, and Li Guo.	
948	2023. Divide, conquer, and combine: Mixture of	
949	semantic-independent experts for zero-shot dialogue	
950	state tracking. <i>arXiv preprint arXiv:2306.00434</i> .	
	Weizhi Wang, Zhirui Zhang, Junliang Guo, Yinpei Dai,	951
	Boxing Chen, and Weihua Luo. 2022. Task-oriented	952
	dialogue system as natural language generation. In	953
	<i>Proceedings of the 45th International ACM SIGIR</i>	954
	<i>Conference on Research and Development in Infor-</i>	955
	<i>mation Retrieval</i> , pages 2698–2703.	956
	Tsung-Hsien Wen, David Vandyke, Nikola Mrkšić, Mil-	957
	ica Gašić, Lina M. Rojas-Barahona, Pei-Hao Su, Ste-	958
	fan Ultes, and Steve Young. 2017. A network-based	959
	end-to-end trainable task-oriented dialogue system .	960
	In <i>Proceedings of the 15th Conference of the Euro-</i>	961
	<i>pean Chapter of the Association for Computational</i>	962
	<i>Linguistics: Volume 1, Long Papers</i> , pages 438–449,	963
	Valencia, Spain. Association for Computational Lin-	964
	guistics.	965
	Jason D Williams, Kavosh Asadi, and Geoffrey	966
	Zweig. 2017. Hybrid code networks: practical	967
	and efficient end-to-end dialog control with super-	968
	vised and reinforcement learning. <i>arXiv preprint</i>	969
	<i>arXiv:1702.03274</i> .	970
	Chien-Sheng Wu, Steven C.H. Hoi, Richard Socher,	971
	and Caiming Xiong. 2020. TOD-BERT: Pre-trained	972
	natural language understanding for task-oriented di-	973
	alogue . In <i>Proceedings of the 2020 Conference on</i>	974
	<i>Empirical Methods in Natural Language Processing</i>	975
	<i>(EMNLP)</i> , pages 917–929, Online. Association for	976
	Computational Linguistics.	977
	Chien-Sheng Wu, Andrea Madotto, Ehsan Hosseini-	978
	Asl, Caiming Xiong, Richard Socher, and Pascale	979
	Fung. 2019. Transferable multi-domain state genera-	980
	tor for task-oriented dialogue systems. <i>arXiv preprint</i>	981
	<i>arXiv:1905.08743</i> .	982
	Qingyang Wu, James Gung, Raphael Shu, and Yi Zhang.	983
	2023. Diacttod: Learning generalizable latent di-	984
	alogue acts for controllable task-oriented dialogue	985
	systems . <i>arXiv preprint arXiv:2308.00878</i> .	986
	Kun Qian Zhou Yu Xiao Yu, Qingyang Wu. 2022. Krls:	987
	Improving end-to-end response generation in task	988
	oriented dialog with reinforced keywords learning .	989
	<i>arXiv preprint arXiv:2211.16773</i> .	990
	Yunyi Yang, Yunhao Li, and Xiaojun Quan. 2021. Ubar:	991
	Towards fully end-to-end task-oriented dialog system	992
	with gpt-2 . In <i>Proceedings of the AAAI Conference</i>	993
	<i>on Artificial Intelligence</i> , volume 35, pages 14230–	994
	14238.	995
	Xiaoxue Zang, Abhinav Rastogi, Srinivas Sunkara,	996
	Raghav Gupta, Jianguo Zhang, and Jindong Chen.	997
	2020. Multiwoz 2.2: A dialogue dataset with addi-	998
	tional annotation corrections and state tracking base-	999
	lines. In <i>Proceedings of the 2nd Workshop on Natu-</i>	1000
	<i>ral Language Processing for Conversational AI, ACL</i>	1001
	<i>2020</i> , pages 109–117.	1002
	Yizhe Zhang, Siqi Sun, Michel Galley, Yen-Chun Chen,	1003
	Chris Brockett, Xiang Gao, Jianfeng Gao, Jingjing	1004
	Liu, and Bill Dolan. 2020. DIALOGPT : Large-scale	1005
	generative pre-training for conversational response	1006

1007 [generation](#). In *Proceedings of the 58th Annual Meet-*
1008 *ing of the Association for Computational Linguistics:*
1009 *System Demonstrations*, pages 270–278, Online. As-
1010 sociation for Computational Linguistics.

1011 Jeffrey Zhao, Raghav Gupta, Yuan Cao, Dian Yu,
1012 Mingqiu Wang, Harrison Lee, Abhinav Rastogi,
1013 Izhak Shafran, and Yonghui Wu. 2022. Description-
1014 driven task-oriented dialog modeling. *arXiv preprint*
1015 *arXiv:2201.08904*.

1016 Yanqi Zhou, Tao Lei, Hanxiao Liu, Nan Du, Yanping
1017 Huang, Vincent Y Zhao, Andrew M. Dai, Zhifeng
1018 Chen, Quoc V Le, and James Laudon. 2022. [Mixture-](#)
1019 [of-experts with expert choice routing](#). In *Advances*
1020 *in Neural Information Processing Systems*.

1021 **A Statistics of Results**

Model	Dataset	Module	Metric	Mean	Std
T5_{small}	MultiWOZ 2.1	NLU	JGA	59.69	0.028
		NLG	Inform	92.50	0.167
			Sucess	74.00	0.335
			BLEU	16.89	0.019
		Combined	100.14	-	
	MultiWOZ 2.2	NLU	JGA	59.60	0.026
		NLG	Inform	89.6	0.207
			Sucess	76.2	0.349
BLEU			17.1	0.031	
	Combined	100.1	-		
T5_{base}	MultiWOZ 2.1	NLU	JGA	60.36	0.017
		NLG	Inform	92.3	0.071
			Sucess	78.8	0.217
			BLEU	16.88	0.011
		Combined	102.43	-	
	MultiWOZ 2.2	NLU	JGA	60.08	0.026
		NLG	Inform	89.0	0.182
			Sucess	76.0	0.349
BLEU			17.6	0.013	
	Combined	99.7	-		

Table 5: Mean and standard deviation of all reported scores in Table 2 and Table 3 using 5 randomly sampled batch sizes, which are the same for all models and datasets. Student paired t-test shows $p < 0.01$ for scores changed by batch size. Combined = (Inform + Success) \times 0.5 + BLEU.