

```

1 %matplotlib inline
2 import pandas as pd
3 import numpy as np
4 import matplotlib
5 from matplotlib import pyplot as plt
6 import string
7 from collections import defaultdict

1 plt.rcParams["font.size"] = 16
2 SHOW_RG = False
3 bao_run_path = f'/Users/saharshbarve/Work/uiuc/study/Semester_3/cs511/project-bao/breakingBao/exp2-query-selectivity/bao_run.t
4 pg_run_path = f'/Users/saharshbarve/Work/uiuc/study/Semester_3/cs511/project-bao/breakingBao/exp2-query-selectivity/pg_run.txt

1 with open(pg_run_path) as f:
2     data = f.read().split("\n")[2:]
3 data = [x.split(" ") for x in data if len(x) > 1 and (x[0] in string.digits or x[0] == "x")]
4
5 data = [(x[0], x[1], float(x[2]), x[3], float(x[4])) for x in data]
6 pg_data = data
7 pg_times = np.array([x[2] for x in pg_data])
8 pg_times -= np.min(pg_times)
9 pg_times /= 60
10
11
12 def read_bao_data(fp):
13     with open(fp) as f:
14         data = f.read().split("\n")[2:]
15
16     training_times = []
17     for idx in range(len(data)):
18         if data[idx].strip().startswith("Initial input channels"):
19             prev_line = data[idx-1].split(" ")
20             if prev_line[0] == "Retry":
21                 continue
22             training_times.append(float(prev_line[2]))
23
24
25     training_times = np.array(training_times)
26
27     data = [x.split(" ") for x in data if len(x) > 1 and (x[0] in string.digits or x[0] == "x")]
28     data = [(x[0], x[1], float(x[2]), x[3], float(x[4])) for x in data]
29     bao_data = data
30
31     bao_times = np.array([x[2] for x in bao_data])
32     training_times -= np.min(bao_times)
33     bao_times -= np.min(bao_times)
34
35     bao_times /= 60
36     training_times /= 60
37     return bao_data, bao_times, training_times
38
39 bao_data, bao_times, training_times = read_bao_data(bao_run_path)
40 if SHOW_RG:
41     bao_rb_data, bao_rb_times, training_rb_times = read_bao_data("bao_with_regblock.txt")

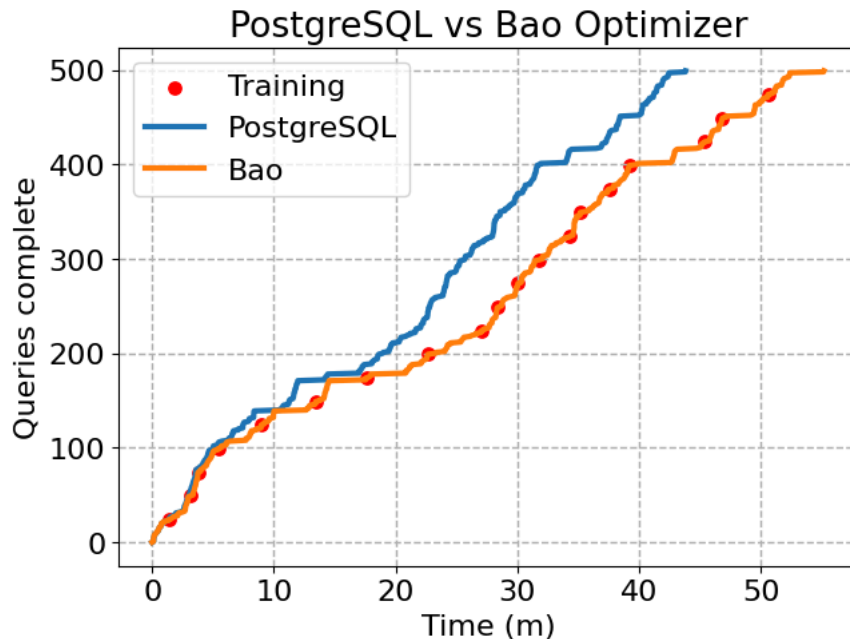
1 queries_complete = np.arange(0, len(pg_times))
2
3 fig, ax = plt.subplots(1, 1, constrained_layout=True)
4
5
6 train_y = []
7 train_rb_y = []
8 for tt in training_times:
9     idx = np.searchsorted(bao_times, tt)
10    train_y.append(idx)
11
12 if SHOW_RG:
13     for tt in training_rb_times:
14         idx = np.searchsorted(bao_rb_times, tt)
15         train_rb_y.append(idx)
16
17 plt.scatter(training_times, train_y, s=45, color="red", label="Training")
18
19 ax.plot(pg_times, queries_complete, label="PostgreSQL", lw=3)

```

```

20 ax.plot(bao_times, queries_complete, label="Bao", lw=3)
21
22 if SHOW_RG:
23     plt.scatter(training_rb_times, train_rb_y, s=45, color="red")
24     ax.plot(bao_rb_times, queries_complete, label="Bao (w/ exploration)", lw=3)
25
26 ax.set_xlabel("Time (m)")
27 ax.set_ylabel("Queries complete")
28 ax.set_title("PostgreSQL vs Bao Optimizer")
29
30 ax.grid(linestyle="--", linewidth=1)
31 ax.legend()
32 fig.savefig("queries_vs_time.svg")

```



```

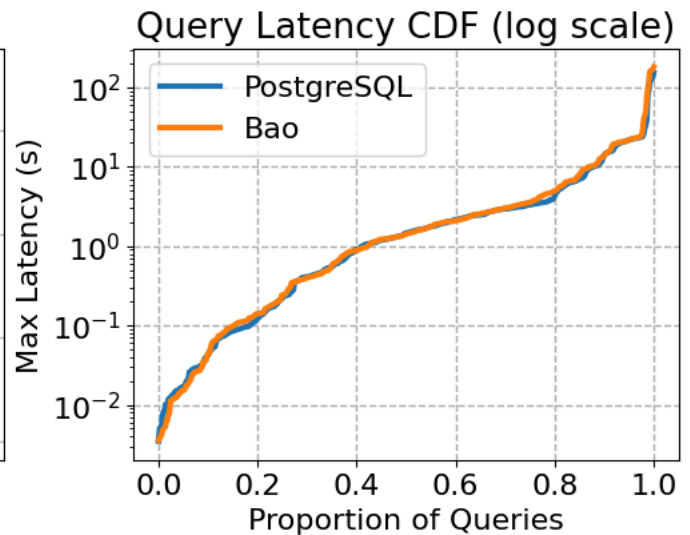
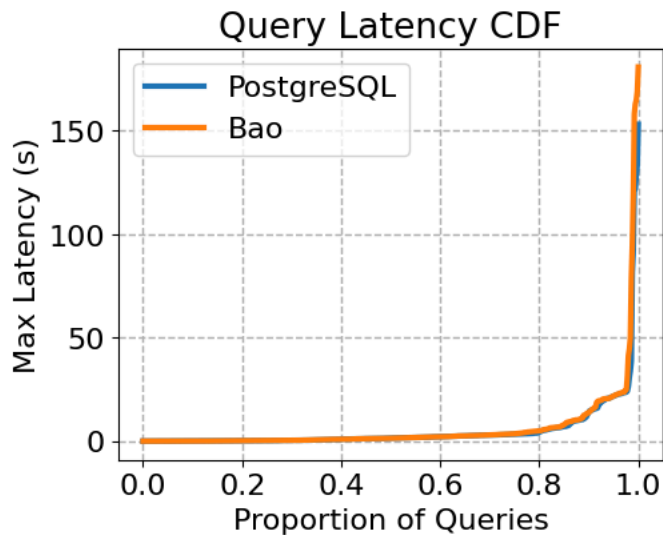
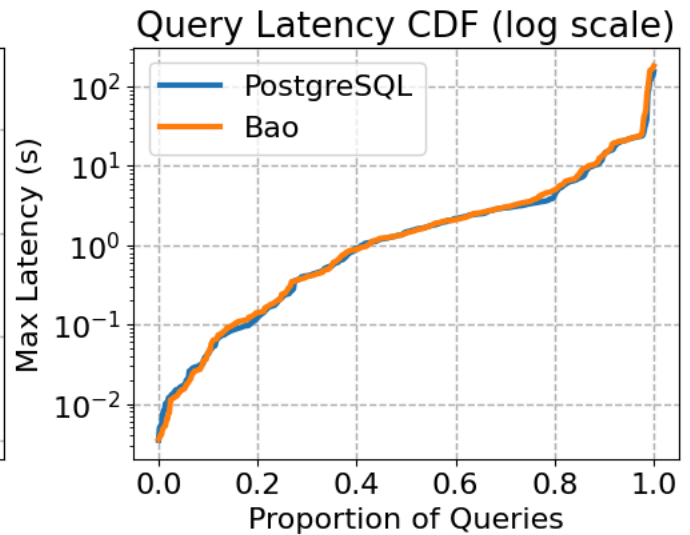
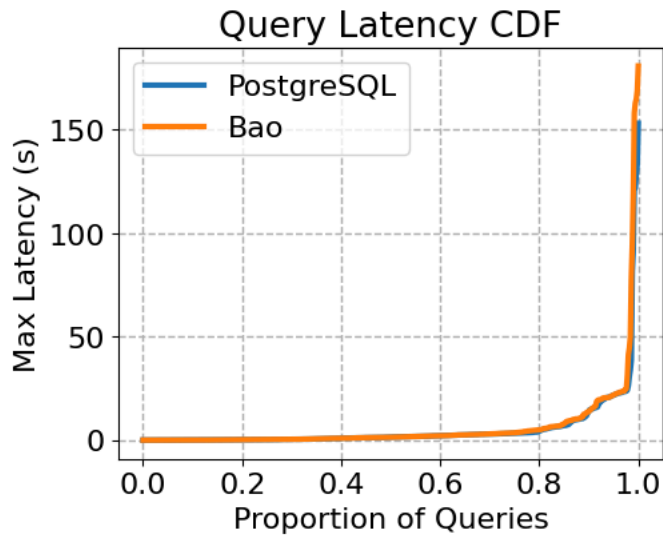
1 all_pg_times = sorted([x[4] for x in pg_data])
2 all_bao_times = sorted([x[4] for x in bao_data])
3
4 if SHOW_RG:
5     all_bao_rb_times = sorted([x[4] for x in bao_rb_data])
6
7
8 fig, axes = plt.subplots(1, 2, figsize=(10, 4), constrained_layout=True)
9
10 ax = axes[0]
11 ax.plot(np.linspace(0, 1, len(all_pg_times)), all_pg_times, lw=3, label="PostgreSQL")
12 ax.plot(np.linspace(0, 1, len(all_pg_times)), all_bao_times, lw=3, label="Bao")
13
14 if SHOW_RG:
15     ax.plot(np.linspace(0, 1, len(all_pg_times)), all_bao_rb_times, lw=3, label="Bao (w/ exploration)")
16
17 ax.grid(linestyle="--", linewidth=1)
18 ax.set_xlabel("Proportion of Queries")
19 ax.set_ylabel("Max Latency (s)")
20 ax.set_title("Query Latency CDF")
21 ax.legend()
22 #ax.set_yscale("log")
23
24
25 ax = axes[1]
26 ax.plot(np.linspace(0, 1, len(all_pg_times)), all_pg_times, lw=3, label="PostgreSQL")
27 ax.plot(np.linspace(0, 1, len(all_pg_times)), all_bao_times, lw=3, label="Bao")
28
29 if SHOW_RG:
30     ax.plot(np.linspace(0, 1, len(all_pg_times)), all_bao_rb_times, lw=3, label="Bao (w/ exploration)")
31
32 ax.grid(linestyle="--", linewidth=1)
33 ax.set_xlabel("Proportion of Queries")
34 ax.set_ylabel("Max Latency (s)")
35 ax.set_title("Query Latency CDF (log scale)")
36 ax.legend()

```

```

37 ax.set_yscale("log")
38 fig.savefig("cdf.svg")
39 fig

```



```

1 # get the last PG time for each query
2 pg_query_time = {}
3 for itm in pg_data:
4     pg_query_time[itm[3]] = itm[4]
5
6 # get each Bao time
7 bao_query_times = defaultdict(list)
8 for itm in bao_data[50:]:
9     bao_query_times[itm[3]].append(itm[4])
10
11 if SHOW_RG:
12     # get each Bao time
13     bao_rb_query_times = defaultdict(list)
14     for itm in bao_rb_data[50:]:
15         bao_rb_query_times[itm[3]].append(itm[4])
16
17 max_repeats = max(len(x) for x in bao_query_times.values())
18
19 def extract_q_number(x):
20     return int(x[x.find("/q")+2:x.find("_", x.find("/q"))])
21
22 q_order = sorted(bao_query_times.keys(), key=extract_q_number)
23
24 grid = [bao_query_times[x] for x in q_order]
25
26 if SHOW_RG:
27     grid_rb = [bao_rb_query_times[x] for x in q_order]
28

```

```

29
30 reg_data = []
31 for idx, q in enumerate(q_order):
32     if SHOW_RG:
33         reg_data.append({"Q": f"q{extract_q_number(q)}",
34                         "PG": pg_query_time[q],
35                         "Bao worst": max(grid[idx]),
36                         "Bao best": min(grid[idx]),
37                         "Bao + E worst": max(grid_rb[idx]),
38                         "Bao + E best": min(grid_rb[idx])})
39     else:
40         reg_data.append({"Q": f"q{extract_q_number(q)}",
41                         "PG": pg_query_time[q],
42                         "Bao worst": max(grid[idx]),
43                         "Bao best": min(grid[idx])})
44
45
46
47 def color_regression(col):
48     def c_for_diff(diff):
49         if diff < 2 and diff > -2:
50             return "background-color: white"
51         elif diff > 0.5:
52             return "background-color: #f27281"
53         else:
54             return "background-color: #9ee3ad"
55
56     to_r = [""]
57
58     if SHOW_RG:
59         pg, bao_worst, bao_best, bao_rg_worst, bao_rg_best = col
60     else:
61         pg, bao_worst, bao_best = col
62
63
64     to_r.append(c_for_diff(bao_worst - pg))
65     to_r.append(c_for_diff(bao_best - pg))
66
67     if SHOW_RG:
68         to_r.append(c_for_diff(bao_rg_worst - pg))
69         to_r.append(c_for_diff(bao_rg_best - pg))
70
71     return to_r
72
73 reg_data = pd.DataFrame(reg_data).set_index("Q")
74 reg_data.style.apply(color_regression, axis=1)

```



```

-----
ValueError                                Traceback (most recent call last)
Cell In[7], line 22
    19 def extract_q_number(x):
    20     return int(x[x.find("/q")+2:x.find("_", x.find("/q"))])
--> 22 q_order = sorted(bao_query_times.keys(), key=extract_q_number)
    24 grid = [bao_query_times[x] for x in q_order]
    26 if SHOW_RG:

Cell In[7], line 20, in extract_q_number(x)
    19 def extract_q_number(x):
--> 20     return int(x[x.find("/q")+2:x.find("_", x.find("/q"))])

ValueError: invalid literal for int() with base 10: 'uery'

```

Double-click (or enter) to edit

