

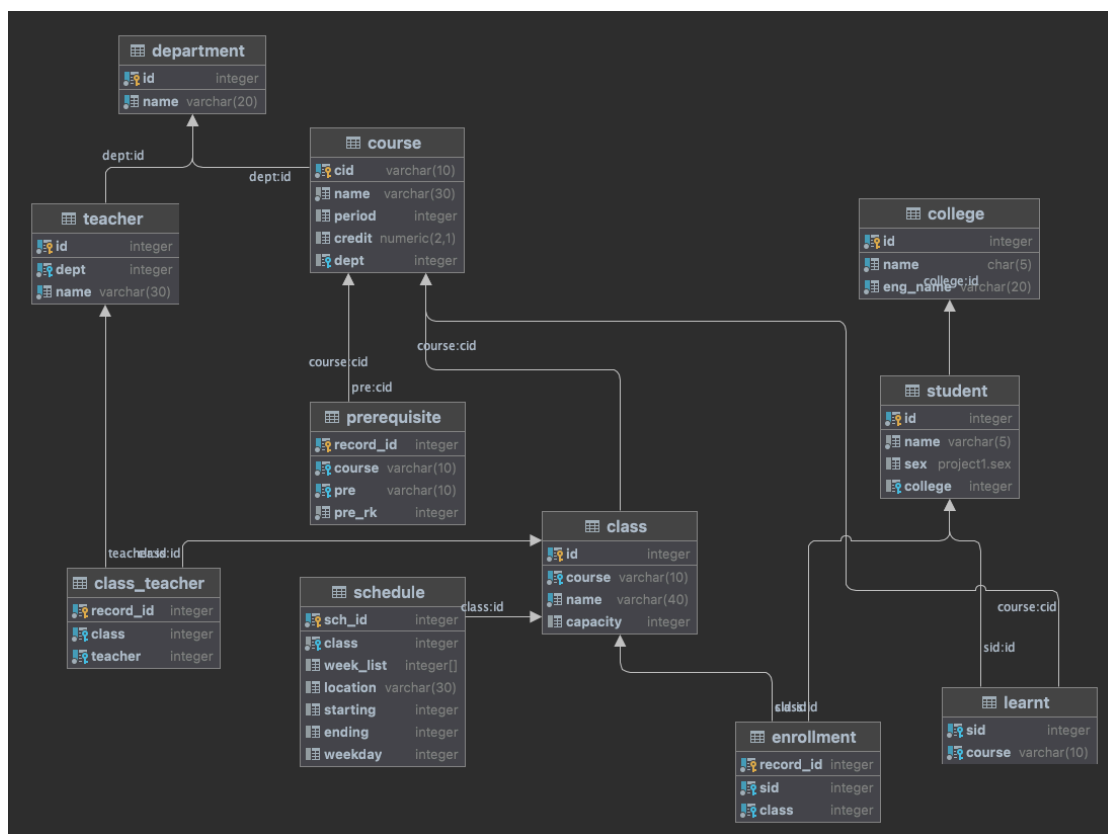
Project 1

Contents

1 结构设计	1
1.1 Department	1
1.2 Teacher	1
1.3 Course	1
1.4 Class	2
1.5 Class-Teacher	2
1.6 Schedule	2
1.7 Prerequisite	3
1.8 College	3
1.9 Student	3
1.10 Learnt	4
1.11 Enrollment	4
2 导入数据	4
2.1 数据清洗	4
2.2 导入课程信息	6
2.3 导入选课信息	11
2.3.1 关闭自动提交	11
2.3.2 批处理及分页大小	12
2.3.3 延迟检查外键约束	13
2.3.4 异步插入	13
2.3.5 连接池	14
3 DML 性能分析	15
3.1 插入条目 (增)	15
3.2 删除条目 (删)	16
3.3 查询条目 (查)	17
3.4 修改条目 (改)	18
4 DBMS-文件系统性能对比	19
4.1 大数据量	19
4.2 小数据量	20

5 附加项	21
5.1 高并发及事务管理	21
5.1.1 高并发	21
5.1.2 事务管理	22
5.2 用户权限管理	25
5.2.1 角色 (Roles)	25
5.2.2 权限授予	25
5.2.3 撤回权限	29
5.2.4 组角色	29
5.2.5 文件系统用户权限管理	30
5.3 数据库索引及文件 IO	30
5.3.1 数据库索引	30
5.3.2 文件索引 (数据结构)	32
5.4 跨平台性能对比	33

1 结构设计



1.1 Department

```

1 CREATE TABLE department (
2     id serial PRIMARY KEY,
3     name varchar(20) UNIQUE NOT NULL
4 );
  
```

无名的部门是无意义的，故为其加上非空约束，学校中的部门均不重名，即 `UNIQUE(name)`。同时，此表有助于我们减少数据冗余、更方便进行部门更名或删除等操作，也允许了无成员的部门存在。

1.2 Teacher

```

1 CREATE TABLE teacher (
2     id serial PRIMARY KEY,
3     dept int NOT NULL REFERENCES department (id),
4     name varchar(30) NOT NULL,
5     CONSTRAINT teacher_unk UNIQUE (dept, name)
6 );
  
```

考虑到学校中可能有重名的教师，而将范围缩小到一个院系里，出现重名的概率大大减少，在无明确指示的情况下我们可以认为部门一致、姓名一致即为同一个人。

1.3 Course

```

1 CREATE TABLE course (
  
```

```
2  cid    varchar(10) PRIMARY KEY,
3  name   varchar(30) NOT NULL,
4  period int CHECK (period > 0),
5  credit numeric(2, 1) CHECK (credit >= 0),
6  dept   int REFERENCES department (id)
7  );
```

本项目共设计三个层次：*course*（共享同一教学大纲）→ *class*（一门课程的多个教学班）→ *schedule*（某一教学班的上课安排）。与 *int* 型主键相比，使用 *varchar* 对后续的查询性能影响不大^[1]，却更加方便查询，均衡二者考虑，选用课程号作为主键。并为 *period*（学时）及 *credit*（学分）加上合理的检查，减少数据库在经过多次数据更改后误操作修改为特别异常的可能性。

1.4 Class

```
1  CREATE TABLE class (
2      id        serial PRIMARY KEY,
3      course    varchar(10) REFERENCES course (cid) NOT NULL,
4      name      varchar(40)                                NOT NULL,
5      capacity  int CHECK (capacity > 0)
6  );
```

这里对课程容量进行检查，但允许其置空。另外，存在同一课程的多个班级同名的情况（FIN204 有两个英文班），故不加 *UNIQUE* 约束。

1.5 Class-Teacher

```
1  CREATE TABLE class_teacher (
2      record_id serial PRIMARY KEY,
3      class     int REFERENCES class (id) NOT NULL,
4      teacher   int REFERENCES teacher (id) NOT NULL,
5      CONSTRAINT cls_tec_unk UNIQUE (class, teacher)
6  );
```

检查数据我们发现，部分班级缺失教师信息，但有教学班同时有数个教师。为了遵循第一范式，我们需要为某班级的所有教师均插入一条记录，且由实际情况可知这不会造成过度的数据膨胀。相比起在 *Class* 表中使用数组存教师，使用此表能更加方便的对教师安排进行调动，且能很方便的查询某门课的所有教师或某教师所教授的所有课程。

1.6 Schedule

```
1  CREATE TABLE schedule (
2      sch_id    serial PRIMARY KEY,
3      class     int REFERENCES class (id) NOT NULL,
4      week_list int[],
5      location  varchar(30),
6      starting  int CHECK (starting >= 1 AND starting <= 12),
7      ending    int CHECK (ending >= 1 AND ending <= 12),
8      weekday   int CHECK (weekday >= 1 AND weekday <= 7),
9      CONSTRAINT time_valid CHECK (ending >= starting)
```

```
10 );
```

我们可以为课程时间进行细致的检查（如上四处 CHECK 所示），另外，考虑到周数列表往往平均长度超过 10，若一味追求 1NF，表中将会存入超过教学班数目十倍的记录，且不利于删除，因此将其保存为数组，这样亦能很方便的检查、调整上课周数。

1.7 Prerequisite

```
1 CREATE TABLE prerequisite (
2     record_id serial PRIMARY KEY,
3     course    varchar(10) REFERENCES course (cid) NOT NULL,
4     pre       varchar(10) REFERENCES course (cid) NOT NULL,
5     pre_rk    int                                NOT NULL,
6     CONSTRAINT preq_unk UNIQUE (course, pre)
7 );
```

课程的先修要求往往是 $(A_1 \vee \dots \vee A_n) \wedge \dots \wedge (Z_1 \vee \dots \vee Z_n)$ 的形式，一般可将其先以“并且”分割为数块以“或者”链接的子关系。此时我们将“或块”分别标记为 *rank*，每个 *rank* 中只需已学 *n* 一门课，且满足所有 *rank* 均学过即可。

1.8 College

```
1 CREATE TABLE college (
2     id        serial PRIMARY KEY,
3     name      char(5) UNIQUE NOT NULL,
4     eng_name  varchar(20) UNIQUE NOT NULL
5 );
```

存学生书院信息，这里将中英文名分别存放，并认为学校中所有书院的中、英文名均不相同。

1.9 Student

```
1 CREATE TYPE sex AS enum ('M', 'F');
2
3 CREATE TABLE student (
4     id        serial PRIMARY KEY,
5     name      varchar(5) NOT NULL,
6     sex       sex,
7     college   int REFERENCES college (id) DEFERRABLE
8 );
```

要保存一条学生记录应该至少知道学生姓名，故将 *name* 条目设为非空约束，而其余条目非关键信息可以为空；为了体现要求中 "Use appropriate types for different fields of data" 且 "be as easy to expand as possible"，选择用自定义枚举类保存性别，这相比起使用 *bool* 标识男女更加直观，且为将来新增性别增加了灵活性。

1.10 Learnt

```
1 CREATE TABLE learnt (  
2     sid    int REFERENCES student (id) DEFERRABLE NOT NULL,  
3     course varchar(10) REFERENCES course (cid) NOT NULL,  
4     CONSTRAINT learn_unk UNIQUE (sid, course)  
5 );
```

数据中未给出选课具体教学班，从常理推测此应该为已学过课程。且学生-课程同时存在才作为一条有意义的记录，故两者均设为 NOT NULL，另外认为修完一门课（无论是一次过还是挂科重修）均只留作一条记录，故加上 UNIQUE。

1.11 Enrollment

```
1 CREATE TABLE enrollment (  
2     record_id serial PRIMARY KEY,  
3     sid    int REFERENCES student (id) DEFERRABLE NOT NULL,  
4     class  int REFERENCES class (id) NOT NULL,  
5     CONSTRAINT enroll_unk UNIQUE (sid, class)  
6 );
```

此表存当前学期选课数据，即在开放选课时经过检查满足先修条件后允许插入，删除一条记录即为退课。在本次项目中此表留空。

2 导入数据

2.1 数据清洗

课程信息 json 文件的原始数据中包含大量空字符等无意义信息，也存在大量不遵循第一范式的条目。首先需要进行数据清洗，并将清洗完的数据存为新 json 文件，方便后面对比文件系统与数据库的性能时可以直接使用。而选课信息格式无需清洗，仅在插入数据库时需要检查是否重复以及有无对应课程。

```
1 import pandas as pd  
2  
3 def clear_data(ori: str, out: str):  
4     with open(ori) as f:  
5         cif = pd.read_json(f)  
6  
7         nif = pd.DataFrame()  
8         nif['totalCapacity'] = cif['totalCapacity']  
9         nif['courseId'] = cif['courseId']  
10        nif['prerequisite'] = [clear_pre(pre) for pre in cif['prerequisite']]  
11        nif['courseHour'] = cif['courseHour']  
12        nif['courseCredit'] = cif['courseCredit']  
13        nif['courseName'] = [clear_str(c) for c in cif['courseName']]  
14        nif['className'] = [clear_str(c) for c in cif['className']]  
15        nif['courseDept'] = cif['courseDept']  
16        nif['teacher'] = [clear_teacher(t) for t in cif['teacher']]  
17        nif['classList'] = clear_class_list(cif['classList'])
```

```

18 with open(out, 'w+', encoding='utf-8') as f:
19     f.write(nif.to_json(orient='records', force_ascii=False))
20

```

我们使用 Pandas 来更灵活的管理数据框架，检查数据后发现需要清洗的数据为以下几类：

courseName 首位存在空字符，且全半角使用不统一。我们将数据统一使用半角符号，查询时也将用户输入的查询条件清洗为半角即可。

```

1 def clear_str(s):
2     if s is None:
3         return None
4     return s.strip() \
5         .replace(' (', '(') \
6         .replace(' ) ', ')') \
7         .replace(', ', ',')

```

```

1 >>> clear_str('文物精品研究鉴赏 ')
2     '文物精品研究鉴赏'
3 >>> clear_str('高等数学 (下) A')
4     '高等数学(下)A'

```

teacher 教师列表不仅存在空字符，也不满足第一范式，我们直接将其处理为 Python 列表。

```

1 def clear_teacher(t):
2     if t is None:
3         return None
4     return [t.strip() for t in clear_str(t).split(',')]

```

```

1 >>> clear_teacher('\t唐珂 , 汤小菊\n , Alan Turing , 于仕琪 ')
2     ['唐珂', '汤小菊', 'Alan Turing', '于仕琪']

```

prerequisite 首先按“与”关系分割为数个“或”关系簇，再将其分割为单个课程。有些课程存在多种形式，并存在重复的多项，但为了体现 DBMS 的性能，这里不做处理，由 DBMS 负责。其难点在于需要清楚逻辑关系的括号，而不能影响到课程名中的括号，否则课程名中带括号或原有括号被删将影响数据库对课程号的匹配。

```

1 def clear_pre(pre):
2     if pre is None:
3         return None
4     return [[ors.strip('() ') + '(' if ors.strip('() ').endswith(('上', '下')) else ''
5             for ors in ands.split('或者')]
6             for ands in pre.split('并且')]

```

```

1 >>> clear_pre('(高等数学 (下) A 或者 高等数学 (下) 或者 数学分析II) 并且 (大学物理A(下) 或者 大
2     学物理 B(下) 或者 大学物理A (下) ) 并且 (线性代数I-A 或者 线性代数I)')
3     [['高等数学(下)A', '高等数学(下)', '数学分析II'], ['大学物理A(下)', '大学物理 B(下)', '大学
4         物理A(下)'], ['线性代数I-A', '线性代数I']]

```

classList 其内部又分为多个可能需要清洗的项：

```

1 def clear_class_list(df):
2     class_list = []
3     for cls in df:
4         cls_items = []
5         for item in cls:
6             item_cls = {}
7             item_cls.update({'weekList': [int(x) for x in item.get('weekList')]})
8             item_cls.update({'location': clear_str(item.get('location'))})
9             item_cls.update({'classTime': item.get('classTime')})
10            item_cls.update({'weekday': item.get('weekday')})
11            cls_items.append(item_cls)
12        class_list.append(cls_items)
13    return class_list

```

The image displays two JSON snippets side-by-side, illustrating the state of course data before and after a cleaning process. The left snippet represents the 'before' state, showing a course entry with a 'weekList' containing integers and a 'location' field with a full, unprocessed address string. The right snippet represents the 'after' state, where the 'weekList' is a list of integers and the 'location' field has been cleaned to a simplified string format.

数据清洗前后对比

2.2 导入课程信息

本部分数据量远小于选课数据，因此我们注重于导入的准确性（优化速度的方案于下节展开讨论）。

由于此部分外键限制错综复杂，若只导一次数据，各种子查询及外键冲突处理等会使得 DML 变得异常复杂且可读性较低，因此，以下我们多次遍历数据，将存在外键约束的条目分别开来。

```

1 import psycopg2 as psql
2 from psycopg2.extras import execute_batch
3 # import other packages
4
5 class DB:
6     def __init__(self, conf: str = './user.yml', data='./data/clear_course_info.json'):
7         with open(conf, 'r', encoding='utf-8') as cfg:
8             conf = yaml.safe_load(cfg)
9         try:
10             self.conn = psql.connect(host=conf['host'], port=conf['port'],
11                                     user=conf['user'], password=conf['pwd'],

```



```

12         database=conf['db'])
13         self.cur = self.conn.cursor()
14         self.conn.autocommit = False
15     except Exception as e:
16         print(e)
17         sys.exit(1)
18     with open(data, 'r') as dat:
19         self.cdata = pd.read_json(dat)
20         print('Database connected && Data loaded to RAM.')
21
22     def __del__(self):
23         self.conn.commit()
24         self.conn.close()

```

这里定义了 DB 类，并在实例化对象时自动连接数据库（连接所需密码、数据库等信息由配置文件读取而得）并将数据载入内存；程序结束销毁对象时进行一次提交并关闭链接。

```

1     def submitter(self, time_start, sql, data):
2         print(f'Preparing data in py ({time.perf_counter() - time_start:.4f}s)', end=' >>> ')
3         try:
4             tm = time.perf_counter()
5             execute_batch(self.cur, sql, data)
6             self.conn.commit()
7             print(f'Submitted {len(data)} requests '
8                   f'({time.perf_counter() - tm:.4f}s | avg={len(data) / (time.perf_counter() - tm)
9                     :.4f}i/s)')
10        except Exception as e:
11            print(e)

```

通用提交插入命令的函数，完成了计时，使用批处理（细节于下节讨论）。

```

1     def create_tables(self, ddl_path):
2         with open(ddl_path, 'r') as dl:
3             sql_list = dl.read().split(';')[::-1]
4         try:
5
6             for sql_item in sql_list:
7                 self.cur.execute(sql_item)
8                 self.conn.commit()
9         except Exception as e:
10            print(repr(e))

```

为了实现高度的自动化，这里将第一节的 DDL 及以下两句命令写入 SQL 文件中。利用 PostgreSQL 命令由分号结束的性质，此函数逐句执行 SQL 文件中的命令（删除原有 schema 并重新建表）。

```

1     DROP SCHEMA IF EXISTS project1 CASCADE;
2     CREATE SCHEMA project1;

```

```

1     @timer
2     def ins_dept(self):

```

```
3     sql = '''
4         INSERT INTO project1.department (name)
5         VALUES (%s)
6         ON CONFLICT DO NOTHING;
7     '''
8
9     tm = time.perf_counter()
10    dept = [(x,) for x in self.cdata['courseDept']]
11    self.submitter(tm, sql, dept)
12
13    @timer
14    def ins_teacher(self):
15        sql = '''
16            INSERT INTO project1.teacher(name, dept)
17            VALUES (%s, (SELECT id
18                          FROM project1.department
19                          WHERE name = %s))
20            ON CONFLICT DO NOTHING;
21        '''
22
23        tm = time.perf_counter()
24        teacher = []
25        for index, item in self.cdata.iterrows():
26            if item['teacher'] is None:
27                continue
28            for t in item['teacher']:
29                teacher.append((t, item['courseDept']))
30        self.submitter(tm, sql, teacher)
31
32    @timer
33    def ins_course(self):
34        sql = '''
35            INSERT INTO project1.course
36            VALUES (%s, %s, %s, %s, (SELECT id
37                                      FROM project1.department
38                                      WHERE name = %s))
39            ON CONFLICT DO NOTHING;
40        '''
41
42        tm = time.perf_counter()
43        course = [(item['courseId'], item['courseName'], item['courseHour'],
44                    item['courseCredit'], item['courseDept'])
45                  for index, item in self.cdata.iterrows()]
46        self.submitter(tm, sql, course)
47
48    @timer
49    def ins_pre(self):
50        sql = '''
51            WITH preid AS (SELECT cid
52                          FROM project1.course
53                          WHERE name = %s)
54            INSERT INTO project1.prerequisite (course, pre, pre_rk)
55            (SELECT %s, (SELECT cid FROM preid LIMIT 1), %s
```

```

53         WHERE (SELECT COUNT(cid) FROM preid) > 0)
54     ON CONFLICT DO NOTHING;
55     '''
56     tm = time.perf_counter()
57     pres = []
58     for index, item in self.cdata.iterrows():
59         if item['prerequisite'] is None:
60             continue
61         for rk, ands in enumerate(item['prerequisite'], start=1):
62             for ors in ands:
63                 pres.append((ors, item['courseId'], rk))
64     self.submitter(tm, sql, pres)

```

实现避免插入重复一般有两种方案：① INSERT INTO <table> (SELECT <values> WHERE (SELECT COUNT(*) FROM <table>) = 0, 可见该方法不仅书写复杂，本地为了填充批处理命令也需要准备更多数据；② INSERT INTO <table> VALUES (<values>) ON CONFLICT DO NOTHING, 不仅书写简单，填充数据少，且 DBMS 计划任务较简单，执行更快。

id	name
1	体育中心
2	35 电子与电气工程系
3	37 化学系
4	38 生物系
5	40 环境科学与工程学院
6	43 计算机科学与工程系
7	44 机械与能源工程系
8	54 材料科学与工程系
9	56 医学院
10	58 思想政治教育与理论
11	67 地球与空间科学系
12	68 海洋科学与工程系
13	69 数学系
14	79 生物医学工程系
15	81 力学与航空航天工程
16	91 物理系
17	112 金融系
18	410 社会科学中心
19	412 语言中心
20	416 艺术中心
21	420 人文科学学院
22	427 创新创业学院
23	448 高等教育研究中心

id	dept	name
1	1	1 体育部
2	2	1 数学
3	3	1 物理系
4	4	1 体育与健康教育
5	6	1 历史
6	7	1 化学系
7	9	1 生物
8	13	1 计算机
9	17	1 医学
10	19	1 物理
11	22	1 数学
12	35	1 电子与电气工程系
13	37	1 化学系
14	38	1 生物系
15	39	1 环境科学与工程系
16	40	1 环境科学与工程系
17	41	1 计算机科学与工程系
18	42	1 机械与能源工程系
19	43	1 计算机科学与工程系
20	44	1 机械与能源工程系
21	45	1 计算机科学与工程系
22	46	1 计算机科学与工程系
23	47	1 计算机科学与工程系
24	48	1 计算机科学与工程系

id	name	period	credit	dept
1	EE232	32	1.0	1
2	EE332	44	3.0	35
3	CH215	48	3.0	37
4	BI0320	48	3.0	38
5	ESE212	32	2.0	40
6	CH005-2	44	2.0	37
7	EE218	48	3.0	35
8	CS402	16	1.0	43
9	ME306	44	3.0	44
10	ESE210	16	0.5	40
11	CH333	32	2.0	37
12	ESE216	48	3.0	40
13	BI0340	48	3.0	38
14	ME303	48	3.0	44
15	CS401	44	3.0	43
16	MSE202	48	3.0	54
17	ME307	48	3.0	56
18	EE210	32	2.0	35
19	EE330	32	2.0	35
20	ESE302	48	3.0	47
21	OCE366	48	3.0	48
22	MA314	48	3.0	69
23	MA2010	80	4.0	69
24	CH005-1	32	2.0	37

record_id	course	dept	prereq_id	prereq_course
1	EE332	EE202-17	1	EE332
2	7 BI0320	BI0102A	1	7 BI0320
3	11 ESE212	CH101A	1	11 ESE212
4	12 ESE212	CH101B	1	12 ESE212
5	13 ESE212	PHY105A	1	13 ESE212
6	14 ESE212	PHY105B	1	14 ESE212
7	16 ESE212	CH102-17	1	16 ESE212
8	17 CH305-2	CH305-1	1	17 CH305-2
9	18 ME306	ME303	1	18 ME306
10	19 ME306	ME307	1	19 ME306
11	21 ESE218	ESE216	1	21 ESE218
12	22 CH215	CH206	1	22 CH215
13	23 ESE216	MA102B	1	23 ESE216
14	24 ESE216	MA102A	1	24 ESE216
15	25 ESE216	PHY105A	1	25 ESE216
16	26 ESE216	PHY105B	1	26 ESE216
17	28 ME303	MAE202	1	28 ME303
18	31 ME303	ME102	1	31 ME303
19	32 CS401	CS102A	1	32 CS401
20	33 CS401	CS201	1	33 CS401
21	34 CS401	CS202	1	34 CS401
22	44 MSE202	CH101A	1	44 MSE202
23	45 MSE202	MA102B	1	45 MSE202
24	46 MSE202	MA102A	1	46 MSE202

上示正确插入后的数据存储

```

1 @timer
2 def ins_class_te_sc(self):
3     sql_cls = '''
4         INSERT INTO project1.class
5         VALUES (%s, %s, %s, %s);
6     '''
7     sql_tec = '''
8         WITH tid AS (
9         SELECT id
10        FROM project1.teacher

```

```

11         WHERE name = %s
12         AND dept = (SELECT id FROM project1.department WHERE name = %s)
13     )
14     INSERT INTO project1.class_teacher (class, teacher)
15     VALUES (%s, (SELECT id FROM tid));
16     '''
17 sql_sch = '''
18     INSERT INTO project1.schedule (class, week_list, location, starting, ending,
19     weekday)
20     VALUES (%s, %s::int[], %s, %s, %s, %s);
21     '''
22 tm = time.perf_counter()
23 cls = []
24 tec = []
25 sch = []
26 self.cur.execute('''
27     SELECT GREATEST((SELECT last_value - 1
28     FROM project1.class_id_seq),
29     (SELECT MAX(id)
30     FROM project1.class));
31     ''')
32 for idx, (index, item) in enumerate(self.cdata.iterrows(), start=self.cur.fetchone()[0]
33 + 1):
34     cls.append((idx, item['courseId'], item['className'], item['totalCapacity']))
35     if item['teacher'] is not None:
36         for x in iter((t, item['courseDept'], idx) for t in item['teacher']):
37             tec.append(x)
38         for x in iter((idx,
39             str(cl.get('weekList')).replace('[', '{').replace(']', '}'),
40             cl.get('location'),
41             int(cl.get('classTime').split('-')[0]),
42             int(cl.get('classTime').split('-')[1]),
43             cl.get('weekday'))
44             for cl in item['classList']):
45             sch.append(x)
46 print(f'Preparing data in py ({time.perf_counter() - tm:.4f}s)', end=' >>> ')
47 try:
48     tm = time.perf_counter()
49     execute_batch(self.cur, sql_cls, cls)
50     execute_batch(self.cur, sql_tec, tec)
51     execute_batch(self.cur, sql_sch, sch)
52     self.cur.execute('SELECT max(id) FROM project1.class')
53     self.cur.execute('ALTER SEQUENCE project1.class_id_seq RESTART WITH %s;',
54                     self.cur.fetchone()[0] + 1)
55     self.conn.commit()
56     print(f'Submitted {len(cls) + len(tec) + len(sch)} requests '
57           f'({time.perf_counter() - tm:.4f}s | '
58           f'avg={({len(cls) + len(tec) + len(sch)) / (time.perf_counter() - tm):.4f}i/s)'
59           )
60 except Exception as e:

```



```

15         FROM project1.college
16         WHERE name = %s))
17     ON CONFLICT DO NOTHING;''',
18     (stu[3], stu[0], stu[1], stu[2].split('(')[0]))
19     for sel in stu[4:]:
20         self.cur.execute(''INSERT INTO project1.learnt
21             (SELECT %s, %s
22              WHERE EXISTS(SELECT id FROM project1.student WHERE id = %s)
23              AND EXISTS(SELECT cid FROM project1.course WHERE cid = %s
24                          ))
25             ON CONFLICT DO NOTHING;''',
26         (stu[3], sel, stu[3], sel))
27 except Exception as e:
28     print(e)

```

需要注意的是 psycopg2 在建立连接后默认关闭自动提交。实验对比可知相比于 autocommit 时的插入速度，关闭其将性能提升了约 36.7%。

```

Database connected && Data loaded to RAM.
2%|██████████| 80879/3999920 [01:55<1:39:23, 657.16it/s]

```

(a) autocommit = true

```

Database connected && Data loaded to RAM.
100%|██████████| 3999920/3999920 [1:04:15<00:00, 1037.53it/s]
bad_loader cost 3855.280791 sec

```

(b) autocommit = false, bad_loader

分析其原因，自动提交时需要频繁检查数据是否满足 ACID，越早插入的数据受越多次检查，从而耗费大量时间；而关闭自动提交大大减少了检测的次数以及每行被检测的平均次数。

2.3.2 批处理及分页大小

```

1 @timer
2 def batch_loader(self, pg: int):
3     # prepare data
4     try:
5         execute_batch(self.cur,
6             '''INSERT INTO project1.college(name, eng_name)
7             SELECT %s, %s
8             ON CONFLICT DO NOTHING;''',
9             college,
10            page_size=pg)
11     # similar for student info and course selection
12 except Exception as e:
13     print(e)

```

使用批处理 (psycopg2.extras.execute_batch) 使得一次操作中可以执行多条 SQL 语句，相比于一次一次执行效率会提高很多。理论上在一定范围内，batch size (page_size) 越大，批量操作需要向数据库发送请求的次数越少，速度就越快；以下实验测得当分页大小为 1000 左右时相较于较小的分页大小有较好的提

速效果，但当其进一步加大时（下图 (d)），速度反而下降，相较于减少请求次数带来的时间损耗，单次发送请求过大反而降低效果。本项可最多带来约 21.7% 的性能提升。

```
Database connected && Data loaded to RAM.
prepdta: 100%|██████████| 3999920/3999920 [00:10<00:00, 369097.39it/s]
college: 100%|██████████| 3999920/3999920 [04:41<00:00, 14209.64it/s]
student: 100%|██████████| 3999920/3999920 [04:06<00:00, 16224.12it/s]
selectc: 100%|██████████| 14000158/14000158 [28:04<00:00, 8309.01it/s]
batch_loader cost 2225.600311 sec
```

(a) batch size = 100

```
Database connected && Data loaded to RAM.
prepdta: 100%|██████████| 3999920/3999920 [00:10<00:00, 371984.56it/s]
college: 100%|██████████| 3999920/3999920 [04:35<00:00, 14503.55it/s]
student: 100%|██████████| 3999920/3999920 [04:32<00:00, 14662.72it/s]
selectc: 100%|██████████| 14000158/14000158 [27:17<00:00, 8548.39it/s]
batch_loader cost 2198.953124 sec
```

(b) batch size = 500

```
Database connected && Data loaded to RAM.
prepdta: 100%|██████████| 3999920/3999920 [00:11<00:00, 339267.11it/s]
college: 100%|██████████| 3999920/3999920 [04:49<00:00, 13839.47it/s]
student: 100%|██████████| 3999920/3999920 [03:44<00:00, 17851.43it/s]
selectc: 100%|██████████| 14000158/14000158 [20:16<00:00, 11504.03it/s]
batch_loader cost 1743.906198 sec
```

(c) batch size = 1000

```
Database connected && Data loaded to RAM.
prepdta: 100%|██████████| 3999920/3999920 [00:11<00:00, 358783.16it/s]
college: 100%|██████████| 3999920/3999920 [04:39<00:00, 14317.89it/s]
student: 100%|██████████| 3999920/3999920 [04:36<00:00, 14470.30it/s]
selectc: 100%|██████████| 14000158/14000158 [27:35<00:00, 8456.21it/s]
batch_loader cost 2224.802974 sec
```

(d) batch size = 1500

2.3.3 延迟检查外键约束

呼应前面建表中的 *deferrable*，只需要在 *batch* 的基础上加上以下两句就可以在插入事务块中暂时不检查外键约束，而在插入完成后统一检查。

```
1 self.cur.execute('BEGIN TRANSACTION;')
2 self.cur.execute('SET CONSTRAINTS ALL DEFERRED;')
```

如下图所示，尤其是选课数据、学生书院等在插入时能确保满足外键约束的，原来每次插入均花费一定开销检查是否满足外键约束，而延迟检查能节省大量时间。最终用时相比上一步的批处理再次减少 46.15%。

```
Database connected && Data loaded to RAM.
prepdta: 100%|██████████| 3999920/3999920 [00:09<00:00, 403848.22it/s]
college: 100%|██████████| 3999920/3999920 [01:53<00:00, 35293.06it/s]
student: 100%|██████████| 3999920/3999920 [03:49<00:00, 17452.94it/s]
selectc: 100%|██████████| 14000158/14000158 [09:45<00:00, 23929.84it/s]
defer_loader cost 939.021061 sec
```

2.3.4 异步插入

```
1 async def async_loader(self):
2     with open(conf, 'r', encoding='utf-8') as cfg:
3         conf = yaml.safe_load(cfg)
```

```

4     conn = await asyncpg.connect(host=conf['host'], port=conf['port'],
5                                   user=conf['user'], password=conf['pwd'],
6                                   database=conf['db'])
7
8     # prepare data
9     await conn.execute('BEGIN TRANSACTION;')
10    await conn.execute('SET CONSTRAINTS ALL DEFERRED;')
11    await conn.executemany(''INSERT INTO project1.college(name, eng_name)
12                            SELECT $1, $2
13                            ON CONFLICT DO NOTHING;'', atqdm(college))
14
15    # similar for student info and course selection
16    await conn.execute('COMMIT;')

```

此时耗时进一步缩短，相比于上一项再次提升约 45.15%，速度达到了平均 42718 条/秒。然而，异步对于本地数据库提升远不如连接云数据库时带来的提升比例大：如果提供一个网络服务，同时维护很多网络连接，并且每个网络连接都在操作数据库，这个异步过程就非常有意义，它会把数据操作过程本来应该阻塞等待的时间用来处理其它网络连接。

```

Database connected && Data loaded to RAM.
3999920it [00:11, 343416.67it/s]
100%|██████████| 3999920/3999920 [00:41<00:00, 96894.76it/s]
100%|██████████| 3999920/3999920 [01:06<00:00, 60563.41it/s]
100%|██████████| 14000158/14000158 [06:48<00:00, 34231.92it/s]

```

2.3.5 连接池

```

1  async def pool_loader(self, conf='./user.yml'):
2      # data
3      seq = 0
4      for stu in tqdm(self.cdata):
5          if seq == 500000:
6              seq = 0
7              college.append(college_seq.copy())
8              stu_info.append(stu_info_seq)
9              stu_sel.append(stu_sel_seq)
10             college_seq.clear()
11             stu_info_seq.clear()
12             stu_sel_seq.clear()
13             # preparing data packages (500000 students/pkg)
14             seq += 1
15             college.append(college_seq)
16             stu_info.append(stu_info_seq)
17             stu_sel.append(stu_sel_seq)
18             self.cur.execute('BEGIN TRANSACTION;')
19             self.cur.execute('SET CONSTRAINTS ALL DEFERRED;')
20
21             conf = yaml.safe_load(open(conf, 'r', encoding='utf-8'))
22             async with asyncpg.create_pool(host=conf['host'], port=conf['port'],
23                                             user=conf['user'], password=conf['pwd'],
24                                             database=conf['db']) as pool:
25                 async with pool.acquire() as cur:

```



```

26     async for seq in atqdm(college):
27         await cur.executemany('INSERT INTO project1.college(name, eng_name)
28                                SELECT $1, $2
29                                ON CONFLICT DO NOTHING;', seq)
30     # similar for student info and course selection

```

在异步插入的基础上应用连接池能进一步提升速度。传统的使用多个连接进行插入会伴随大量的建立、断开连接开销，而上节异步的真正优势体现在高并发或者说并行插入。而连接池实现了连接的复用，正好弥补了上面的缺陷，两者配合，我们将学生数据按 50 万人/份进行打包，对每个包分别进行批处理插入，综合了以上探究的加速方法的所有优势，最终插入 $(3999920 \times 2 + 14000158)$ 条数据总耗时仅 216 秒，平均每秒插入 10.2 万条记录，性能约为异步插入的 2.38 倍。

```

Database connected && Data loaded to RAM.
3999920it [00:12, 315757.94it/s]
100%|██████████| 8/8 [00:41<00:00, 5.19s/it]
100%|██████████| 8/8 [00:52<00:00, 6.50s/it]
100%|██████████| 8/8 [02:03<00:00, 14.57s/it]

```

下表总结了各方法在插入选课数据（14000158 条）时的耗时，由于其数据量较大且 DBMS 行为统一（无较多的 ON CONFLICT DO NOTHING 情况），能较好的代表其性能。表中相对提升指应用每项技术相比于不使用所带来的提升，综合性能是此项速度与传统插入的对比。

	传统	关闭自动提交	批处理	延迟检查	异步插入	连接池
耗时 / s	6086	3855	1216	585	408	123
速度 / $it \cdot s^{-1}$	657.16	1037.53	11504.03	23929.84	34231.92	113822.42
相对提升	-	57.88 %	1008.79 %	108.01 %	43.05 %	232.51 %
综合性能	-	157.88 %	1750 %	3641 %	5209 %	17320 %

3 DML 性能分析

3.1 插入条目（增）

插入条目耗时主要源自以下几点：

条件检查 关系型数据库为了保证 ACID 中的一致性（Consistency），会在插入动作后检查新插入的条目是否满足所有设置的条件（这里主要指 CHECK），显然，每次插入所需检查的项目越多，耗时越长。以下控制其余变量保持一致，演示不断增加 CHECK 的数目，并在 Python 段用逐条插入（开启自动提交）所导致的时间消耗增加。

```

1 CREATE TABLE demo1 (
2     val int,
3     CONSTRAINT demo1_chk1 CHECK (val > 0),
4     CONSTRAINT demo1_chk2 CHECK (val % 2 = 0),
5     CONSTRAINT demo1_chk3 CHECK (val < 10000)
6 );

```

	No constraint	chk1	chk1&2	chk1&2&3
插入 500 条数据 (ms)	3593	3649	3776	3778

外键约束 每插入一条数据都需要在外键所指向的表中对列检查是否有对应条目，此检查又与外键所指向的表中条目数有关：总数据量越大，在其中寻找某一可能存在的外键耗时越多。但所幸 PostgreSQL 会自动为外键所指向的列设置 index，减少了检查耗时。

```

1 CREATE TABLE demo2 (
2     val int,
3     cid varchar(20)      -- case 1: no foreign key constraint
4     REFERENCES demo3(id) -- case 2: demo3 contains 2 items
5     REFERENCES courses(cid) -- case 3: courses contains 338 items
6 );

```

	无外键	Ref 2 items	Ref 338 items
插入 500 条数据 (ms)	3529	3804	3818

主键 / 唯一性约束 此检查耗时与当前此表中的条目数量相关，已有的条目越多，检查插入后是否依然满足唯一性的耗时越长。在批量插入中，越晚插入的数据因为已有数据不断增多而插入速度变慢。另外，主键列由于存在 index，较一般的 unique 约束插入慢，见下文分析。

```

1 CREATE TABLE demo3 (
2     cid varchar(20) UNIQUE
3 );

```

数据量（原有 → 插入后）	0 → 500	501 → 1000	1000 → 1500
插入 500 条数据 (ms)	3409	3614	3738

Index 更新 当对表中的数据进行增加、删除和修改的时候，索引也要动态的维护，因此会降低数据的维护（增删改）速度。这也是我们需要避免在频繁修改数据的表中建立 index 的原因。

```

1 CREATE TABLE demo4 (
2     val1 int,
3     val2 int
4 );
5 CREATE INDEX demo4_idx1 ON demo4 (val1); -- for comparison
6 CREATE INDEX demo4_idx2 ON demo4 (val2); -- for comparison

```

	无 index	2 indexes
插入 500 条数据 (ms)	3746	3898

3.2 删除条目（删）

delete 操作一般使用 WHERE 做限定词，这部分的耗时类似于查询（见下文分析）；另外对于被指向的外键表中的数据，一般需要加上 CASCADE 进行级联删除。

```

1 ALTER TABLE project1.learnt
2     DROP CONSTRAINT learnt_sid_fkey;
3 ALTER TABLE project1.learnt
4     ADD FOREIGN KEY (sid) REFERENCES project1.student (id)
5     ON DELETE CASCADE;
6
7 DELETE FROM project1.student
8 WHERE sex = 'F'
9     AND id BETWEEN 11000020 AND 11080020;

```

此时的删除学生操作会将该学生的学习记录级联删除，执行删除语句的返回结果为 *40,220 rows affected in 1 s 591 ms*。

3.3 查询条目 (查)

子查询 子查询的设计会较大程度上的影响 DBMS 的性能，因为与其余常规查询操作不同，子查询一般无法被优化器综合性的优化进整体查询中。

```

1 SELECT s.id, s.name
2 FROM project1.student s
3 WHERE s.college = (SELECT id FROM project1.college WHERE college.name = '格兰芬多')
4     AND s.sex = 'M';

```

在查询格兰芬多的所有男生信息时，使用子查询每次平均耗时 0.308s，而查询书院编号及按编号查询数据的两个查询总耗时平均 0.221s，这是因为类似于编程中的函数调用栈，当出现子查询时，查询器需要先保存上下文并跳转执行子查询，再将查询结果返回给“调用栈”的上一级查询任务。

QUERY PLAN	
1	Seq Scan on student s (cost=8.17..2145.17 rows=9963 width=14)
2	Filter: ((college = \$0) AND (sex = 'M'::project1.sex))
3	InitPlan 1 (returns \$0)
4	→ Index Scan using college_name_key on college (cost=0.15..8.17 rows=1 width=4)
5	Index Cond: (name = '格兰芬多'::bpchar)

索引 Index 的设计能大大提升查询的性能。不使用 index 时，为了检查 where 子句中的条件，DBMS 必须遍历整个表格中的所有行，并从众多列中选出其需要的部分，而为常用查询的条件部分设置 index 使得 DBMS 只需检查部分字段，大大减少了需要处理的信息量。

```

1 SELECT sch_id
2 FROM project1.schedule
3 WHERE starting = 3
4     AND ending = 4;

```

在执行此查询时，该表上仅有数据库在设置主键列时自动添加的 index，整表一共七个字段，而查询只涉及四个（且无完全包含查询所需字段的 index），故需要提取每行的所有信息，耗时 43ms。下面我们为查询涉及字段添加 index，再次查询仅需 15ms。此时 DBMS 在序列扫描数据时不需要提取全表信息，仅分析了 sch_time_idx 中包含的较少信息。更多关于索引的内容将于下文讨论。

```
1 CREATE INDEX sch_time_idx ON project1.schedule (sch_id, starting, ending);
```

并行执行 并行化语句的执行过程会被规划器分发给多个后台进程去执行。通过这种并行执行方式，PostgreSQL 能充分发挥多核处理器的威力，从而让语句执行更快完成。通过并行化执行所能节省的时间根据数据库载体机的 CPU 核数的多寡而有所不同，在机器性能强大的情况下，并行化所能带来的性能提升可能非常可观。但是在查询任务数据量较少时，强制的并行化执行可能反而延长查询耗时——当每个进程所分配到的子任务规模较小时，创建 worker 的时间开销甚至大于执行查询的时间开销。

通过下面的命令可以强行开启并行化执行模式，对比下例查询，使用 1 worker 时查询用时 1.906s，并行化任务为其分配 2 workers（如下图 explain 所示），耗时 1.308s。

```
1 SET force_parallel_mode = TRUE;
```

```
1 SELECT course, COUNT(*)
2 FROM project1.learnt
3 GROUP BY course;
```

```

QUERY PLAN
1  Finalize GroupAggregate  (cost=162511.64..162576.50 rows=256 width=14)
2    Group Key: course
3    -> Gather Merge  (cost=162511.64..162571.38 rows=512 width=14)
4        Workers Planned: 2
5        -> Sort  (cost=161511.61..161512.26 rows=256 width=14)
6            Sort Key: course
7            -> Partial HashAggregate  (cost=161498.82..161501.38 rows=256 width=14)
8                Group Key: course
9                -> Parallel Seq Scan on learnt  (cost=0.00..132631.88 rows=5773388 width=6)
```

where 条件数 若每层条件筛选后并不会大幅度减少剩余条目数量时，筛选条件越多，查询结果中剩余的项目分摊的执行时间越多，查询越慢。下例中标注的耗时指在原有条件的基础上增加此条件的耗时。

```
1 SELECT sch_id
2 FROM project1.schedule
3 WHERE starting = 3 -- execution: 8ms
4     AND ending = 4  -- execution: 9ms
5     AND weekday = 3; -- execution: 11ms
```

3.4 修改条目 (改)

一方面，关系型数据库要求修改后的数据依然保持所要求的一致性，并可能需要维护 index，此部分检查的开销的分析见上文中“增”部分；另一方面则是 WHERE 限定词的查询开销。

```
1 UPDATE project1.student
2 SET sex='M'
3 WHERE id BETWEEN 11000020 AND 11080020;
4 -- 39,781 rows affected in 162 ms
```

4 DBMS-文件系统性能对比

4.1 大数据量

这里以提供的选课数据 csv 为例，文件系统的查询使用 Python，与数据库类似的，这里不将文件中内容读取至内存中，而是直接移动文件指针对磁盘中的信息进行处理（在 DBMS 中创建表，实际上相当于建立了数据到硬盘的一种映射，从而加速硬盘到内存的读写）。考虑到 DBMS 在插入信息时需要向多表插入，且涉及大量检查，而文件系统的插入仅需对文件追加写入一行记录，删除与修改同理，两者本质差异过大而无较大的对比价值，因此本部分主要讨论数据的查询性能。

DBMS 以其优化器及性能著称。当使用文件系统查询时，或 directly 对文件进行读写，或将其模仿数据库中的结果加载进内存并进行查询，但均需要大量代码实现且需要手动优化。以下面模糊查询所有曾学过与 C++ 相关的任何课程的学生为例：

```
1 @timer
2 def db_query():
3     for _ in range(10):
4         cur.execute('''SELECT s.id, s.name
5                        FROM project1.learnt l
6                        JOIN project1.student s ON l.sid = s.id
7                        JOIN project1.course c ON l.course = c.cid
8                        WHERE c.name ~ 'C\+\+';''')
9     print(len(cur.fetchall()))
10
11 @timer
12 def fs_query():
13     result_set = []
14     for _ in range(10):
15         f.seek(0)
16         result_set.clear()
17         cid = {
18             c['courseId']
19             for index, c in course_info.iterrows()
20             if re.match(r'.*C\+\+', c['courseName'])
21         }
22         for stu in select_course:
23             for c in cid:
24                 if c in stu[4:]:
25                     result_set.append((stu[3], stu[0]))
26             break
27     print(len(result_set))
```

```
109680
db_query cost 5.603938 sec
```

(a) 数据库查询

```
109680
fs_query cost 47.850062 sec
```

(b) 文件系统查询

上图所示进行十次查询并计时的结果，查询结果 109680 项。其中数据库每次查询平均耗时 0.56s，仅为文件系统查询（平均耗时 4.79s）的 11.7%。查看 DBMS 的查询规划可以发现，其用了如 Hash join、Parallel

sequence scan 等技术节省时间，而在代码量相似的实现文件系统查询时仅相当于多个 Sequence scan 的嵌套。事实上，查询越复杂（指必要性上的复杂，而非故意将查询方式写得复杂而难以优化），DBMS 的性能优势越明显（尤其是当能应用包括后面讨论的 Index 技术时），简单的查询可能已经难以进一步优化，性能将与文件系统相近。

QUERY PLAN	
1	Gather (cost=1008.67..160950.62 rows=40994 width=14)
2	Workers Planned: 2
3	-> Nested Loop (cost=8.67..155851.22 rows=17081 width=14)
4	-> Hash Join (cost=8.24..147985.28 rows=17081 width=4)
5	Hash Cond: ((l.course)::text = (c.cid)::text)
6	-> Parallel Seq Scan on learnt l (cost=0.00..132631.88 rows=5773388 width=10)
7	-> Hash (cost=8.23..8.23 rows=1 width=6)
8	-> Seq Scan on course c (cost=0.00..8.23 rows=1 width=6)
9	Filter: ((name)::text ~ 'C\+\+'::text)
10	-> Index Scan using student_pkey on student s (cost=0.43..0.46 rows=1 width=14)
11	Index Cond: (id = l.sid)

4.2 小数据量

以下展示基于题目提供的课程信息 json 文件，文件系统部分使用 Python 中的 Pandas 库处理。

```

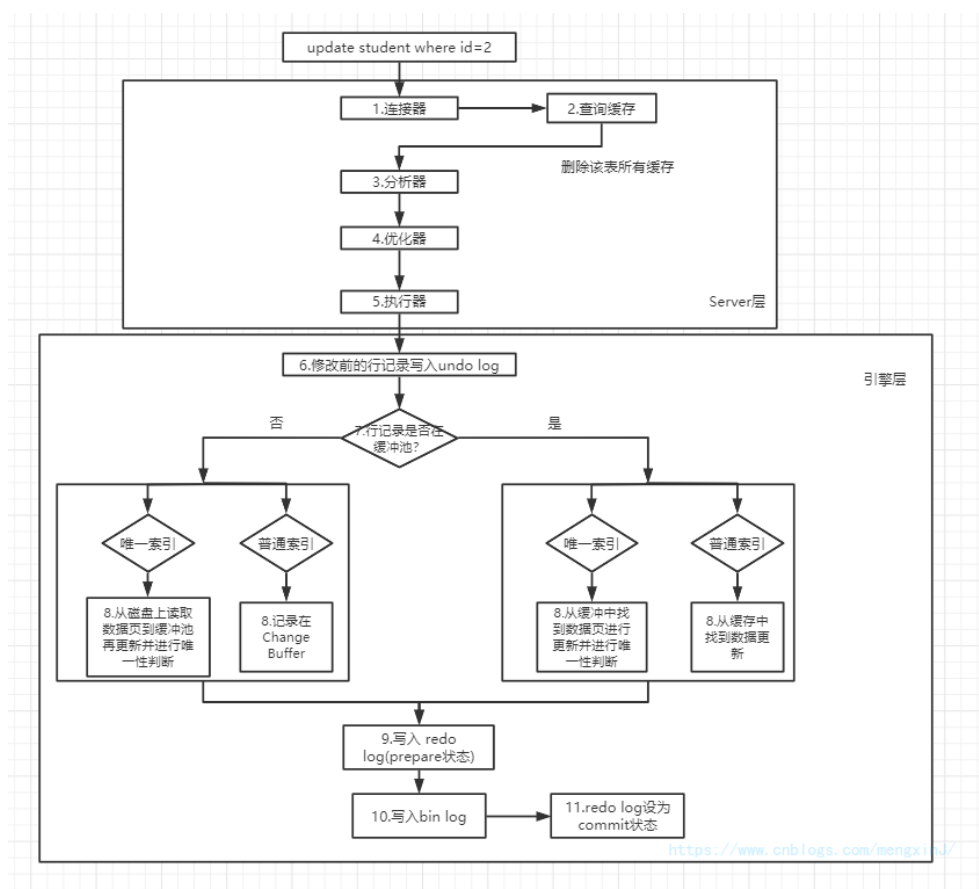
1 @timer
2 def db_query():
3     for _ in range(10):
4         cur.execute('SELECT DISTINCT cid
5                     FROM project1.course
6                     WHERE credit > 2;')
7     print(len(cur.fetchall()))
8
9 @timer
10 def fs_query():
11     result_set = set()
12     for _ in range(10):
13         result_set.clear()
14         for index, c in cif.iterrows():
15             if c['courseCredit'] > 2:
16                 result_set.add(c['courseId'])
17     print(len(result_set))

```

上例重复十次查询所有计算机系开设的课程的列表，数据库耗时 0.298s，文件系统耗时 0.296s，可见即使是在运行速度较慢的 Python 中，依然能以微弱的优势领先数据库。这是由于小文件的遍历开销相对不大，而对于数据库的每次查询都有对任务的优化规划、SQL 的语法检查、用户权限检查等环节的额外开销，因此查询耗时总耗时接近甚至超过文件系统的查询。

数据库的安全性及完善的检查机制是文件系统所远不能及的，但也正是其一系列复杂的流程（语法检查、权限验证）为查询带来了一定的额外开销。当数据量大时，此开销相比于查询开销的比例极小，可以忽略不计，且一系列的优化（哈希、缓存、并行等）大大提升了查询效率；而对于极小的数据量，直接遍历文件已经作为较优的解决方案，DBMS 在规划后可能给出的也是 Seq scan 的方案，但上述额外开销反而

导致耗时较多。但从实际应用中来讲，应用数据库所存储的信息一般至少上千条，且查询的分摊总条目数也远大于遍历依然能跑赢优化的边界数量，也即平均角度上，使用数据库的查询性能是远优于文件系统的；且现实应用离不开数据库带来的一系列完善的检查及权限机制，数据库的优势依然是文件系统无法比拟的。



InnoDB 存储引擎的 SQL 操作过程 [2]

5 附加项

5.1 高并发及事务管理

5.1.1 高并发

支持百万级别的高并发大数据流量已成为现代数据库系统的设计性能标准。在第二部分最后，我们通过异步使用连接池，利用了数据库的高并发支持，在约两分钟的时间段内，数据库持续承受住了 20 个连接的写入操作。关于高并发的一个常见问题是死锁。“加锁（Locking）是数据库在并发访问时保证数据一致性和完整性的主要机制。任何事务都需要获得相应对象上的锁才能访问数据，读取数据的事务通常只需要获得读锁（共享锁），修改数据的事务需要获得写锁（排他锁）。当两个事务互相之间需要等待对方释放获得的资源时，如果系统不进行干预则会一直等待下去，也就是进入了死锁（deadlock）状态。” [3]


```

1 CREATE TABLE _lock (id int, val int);
2 INSERT INTO _lock VALUES (1, 2), (2, 3);
3
4 BEGIN TRANSACTION;
5
6 UPDATE _lock SET val=3 WHERE id = 1;
7
8 UPDATE _lock SET val=5 WHERE id = 2;
9
10
11
12 COMMIT;

1
2
3 BEGIN TRANSACTION;
4
5
6 UPDATE _lock SET val=4 WHERE id = 2;
7
8
9
10 UPDATE _lock SET val=6 WHERE id = 1;
11
12 COMMIT;

```

上图尝试在不同事务中以不同的顺序修改数据，首先在左事务中修改 id 为 1 的数据，此部分被加锁，右事务修改 id 为 2 的数据，同样加锁；第二步让左事务修改 id 为 2 的数据，但此时 DBMS 为避免脏写，已经为其加上写锁，因此左事务的进程需要等待右事务结束以释放写锁；在左事务的等待状态中，右事务修改被左事务添加写锁的数据，此时右事务也进入等待状态；注意到左事务一直在等待右事务结束以释放写锁，而右事务也在等待左事务释放写锁以执行命令并结束，两个事务进入了互相等待对方资源的状态，这种等待是无限的。但现在大部分数据库都能在等待一段时间后检查出死锁，并做出终止一个事务并让另一个事务成功执行的调度（如下图）。但是一个事务长期等待另一个事务释放资源并不会引起 DBMS 的介入。

```

[2021-11-17 23:30:19] [40P01] ERROR: deadlock detected
[2021-11-17 23:30:19] Detail: Process 169114 waits for ShareLock on transaction 123447; blocked by process 168639.
[2021-11-17 23:30:19] Process 168639 waits for ShareLock on transaction 123448; blocked by process 169114.
[2021-11-17 23:30:19] Hint: See server log for query details.
[2021-11-17 23:30:19] Where: while updating tuple (0,1) in relation "_lock"

```

5.1.2 事务管理

事务管理主要为 RDBMS 提供了原子性 (Atomicity)、一致性 (Consistency)、隔离性 (Isolation) 和持久性 (Durability) 这四大特性^[4]。

原子性 指一个事务块中的所有事件或全部成功执行，或全不执行，即在中间任何一个事件失败时需要回滚指事务开始时的状态。在事务块中遇到第一个错误后，Postgres 会忽略后面的命令，直至执行一次 **ROLLBACK** 回滚至事务开始时的状态，或执行 **COMMIT** 并由 DBMS 丢弃此事务块的所有影响。此外，即使所有操作成功执行，用户也可使用 **ROLLBACK** 取消所做修改。

一致性 报告第二部分插入数据即利用了此特性。出于各种原因，人们可能不关心某件事中间过程的高度一致，而仅要求事件前后均满足各种限制即可。“在事务开始之前和事务结束以后，数据库的完整性没有被破坏。这表示写入的资料必须完全符合所有的预设约束、触发器、级联回滚等。”^[5]

隔离性 “数据库允许多个并发事务同时对其数据进行读写和修改的能力，隔离性可以防止多个事务并发执行时由于交叉执行而导致数据的不一致。事务隔离分为不同级别，包括未提交读 (read uncommitted)、提交读 (read committed)、可重复读 (repeatable read) 和串行化 (serializable)。”^[5]

以下测试均在 DataGrip 中使用 console 建立一个连接，并在 CLI 中使用另一角色登录并修改数据，以此模拟并发读写。

读已提交 下面的实验展示在读已提交的隔离级别中，一个查询只能看到查询开始前已提交的数据（查询开始瞬间的数据库快照），无法看到未提交或查询期间由其他事物提交的数据，但查询可看见自身事


```

217 ✓ CREATE TABLE _atomicity (id int);
218
219 ✓ BEGIN TRANSACTION;
220 ✓ INSERT INTO _atomicity VALUES (1);
221 ✗ INSERT INTO _atomicity VALUES ('a');
222 COMMIT;

[22P02] ERROR: invalid input syntax for type integer: "a"
Position: 32

```

(a) 事务块中出现错误

```

[25P02] ERROR: current transaction is aborted, commands ignored until end of transaction block
[22P02] ERROR: invalid input syntax for type integer: "a"
Position: 32

```

(b) 提示当前事务块所有会被取消

```

BEGIN TRANSACTION;
INSERT INTO _atomicity VALUES (1);
INSERT INTO _atomicity VALUES ('a');
COMMIT;
SELECT * FROM _atomicity;

```

(c) 执行 COMMIT，此时 DBMS 取消了事务中的操作，查询结果所示与事务开始前状态一致

```

CREATE TABLE _atomicity (id int);
BEGIN TRANSACTION;
INSERT INTO _atomicity VALUES (1);
SELECT * FROM _atomicity;
ROLLBACK;
SELECT * FROM _atomicity;

```

(d) 左图示成功插入一条数据，执行回滚后再次查询，结果与事务开始前一致

```

CREATE TABLE _consistency (id int UNIQUE DEFERRABLE);
BEGIN TRANSACTION;
SET CONSTRAINTS ALL DEFERRED;
INSERT INTO _consistency VALUES (1);
INSERT INTO _consistency VALUES (1);
COMMIT;

[23505] ERROR: duplicate key value violates unique constraint "_consistency_id_key"
Detail: Key (id)=(1) already exists.

```

(a) 延迟检查使事务块中可违反一些约束，而在 COMMIT 前会检查一致性并报错

务中之前执行的更新。读已提交不会导致脏读（即一个事务读取了另一个未提交的事务写入的数据），同时也是 PostgreSQL 的默认隔离级别。但读已提交也可能造成不可重复读，即同一事务中两个相邻查询返回的结果可能不同，这是因为在两次查询开始的间隔中可能有其他事务被提交。

```

1 conn1=# CREATE TABLE _isolation (id int, val int);
2 conn1=# SELECT * FROM _isolation;
3   id | val
4  ----+-----
5 (0 rows)
6
7 conn1=# BEGIN TRANSACTION;
8 conn1=# SET TRANSACTION ISOLATION LEVEL READ COMMITTED;
9 conn2=# BEGIN TRANSACTION;
10 conn2=# INSERT INTO _isolation VALUES (2, 1);
11 conn1=# SELECT * FROM _isolation;
12   id | val
13  ----+-----

```

```

14 (0 rows)
15
16 conn2=# COMMIT;
17 conn1=# SELECT * FROM _isolation;
18    id | val
19 -----+-----
20     2 |   1
21 (1 row)

```

可重复读 与读已提交的一个非常相似的隔离级别是可重复读，其只允许读取已提交数据，且一个事务两次读取一个数据项期间，不允许其他事务更新该数据，即一个可重复读事务中的查询可以看见在事务中第一个非事务控制语句开始时的一个快照。但其与读已提交一样，依然可能出现幻读。

```

1 conn1=# CREATE TABLE _isolation (id int, val int);
2 conn1=# BEGIN TRANSACTION;
3 conn1=# SET TRANSACTION ISOLATION LEVEL REPEATABLE READ;
4 conn2=# BEGIN TRANSACTION;
5 conn2=# INSERT INTO _isolation VALUES (2, 1);
6 conn2=# COMMIT;
7 conn1=# SELECT * FROM _isolation;
8    id | val
9 -----+-----
10 (0 rows)
11
12 conn1=# COMMIT;
13 conn1=# SELECT * FROM _isolation;
14    id | val
15 -----+-----
16     2 |   1
17 (1 row)

```

可串行化 是四种隔离级别中唯一避免了幻读¹的级别，也即是最高的事务隔离级别，在该级别下，事务串行化顺序执行，可以避免脏读、不可重复读与幻读²。但是这种事务隔离级别效率低下，比较耗数据库性能，不常用。

读未提交 允许读取未提交的数据，是 SQL 允许的最低一致性级别。可能导致脏读、不可重复读以及幻读（但对于 PostgreSQL，此级别的行为与读已提交相同，因为 DBMS 内部实际上只实现了另外三种隔离级别）。另外，以上四种隔离级别均不允许脏写。

持久性 事务块成功结束后，对数据的修改就是永久的，即便系统故障也不应丢失。如下图所示，将 postgres.service 强制突然 kill 掉，模拟系统突然故障，此后重新启动服务，检查数据库存储条目与 kill 前一致（kill 前已结束事务）。

¹一个事务开始后，需要根据数据库中现有的数据做一些更新，于是重新执行一个查询，返回符合查询条件的行，这时发现这些行因为其它最近提交的事务而发生了改变，导致现有事务如果再进行下去可能发发生逻辑上的错误。

²不可重复读对应的是修改，即 UPDATE 操作；而幻读问题对应的是 INSERT 操作。

```
(base) webdepgsngo ➤ sudo systemctl kill postgresql-13.service
(base) webdepgsngo ➤ sudo systemctl status postgresql-13.service
* postgresql-13.service - PostgreSQL 13 database server
   Loaded: loaded (/usr/lib/systemd/system/postgresql-13.service; enabled; vendor preset: disabled)
   Active: inactive (dead) since Wed 2021-11-17 16:11:38 CST; 3s ago
     Docs: https://www.postgresql.org/docs/13/static/
   Process: 1028 ExecStart=/usr/pgsql-13/bin/postmaster -D ${PGDATA} (code=exited, status=0/SUCCESS)
   Process: 993 ExecStartPre=/usr/pgsql-13/bin/postgresql-13-check-db-dir ${PGDATA} (code=exited, status=0/SUCCESS)
   Main PID: 1028 (code=exited, status=0/SUCCESS)

Oct 14 08:27:40 sngo systemd[1]: Starting PostgreSQL 13 database server...
Oct 14 08:27:41 sngo postmaster[1028]: 2021-10-14 08:27:41.522 CST [1028] LOG: redirecting log output to logging collector process
Oct 14 08:27:41 sngo postmaster[1028]: 2021-10-14 08:27:41.522 CST [1028] HINT: Future log output will appear in directory "log".
Oct 14 08:27:42 sngo systemd[1]: Started PostgreSQL 13 database server.
```

5.2 用户权限管理

DBMS 最有用的特性之一便是其完善健全的用户权限管理机制。考虑数据库的日常使用场景，为不同的用户（更正确的说，角色）提供不同级别、不同类型的权限是很有必要的。

5.2.1 角色 (Roles)

PostgreSQL 为管理“用户”及其权限提供了角色机制，在此机制下，“用户”、“用户组”是同一个概念，这点与 MySQL 等数据库的机制有十分关键的区别。

用户与用户组实际上都作为“角色”存储在 `pg_roles` 数据库中，自然，他们拥有同样数目的属性。使用 DBMS 的用户拥有一个登录角色 (*login role*)，同时，它可以从另一个角色继承权限，成为其成员角色 (*member role*)，同时，拥有一些成员角色的角色现在被命名为组角色 (*group role*)。一般来说，出于安全性考虑，组角色不授予登录权限^[7]，因为组角色被设计成一个权限集合，以便于管理多个用户的权限，而不是作为一个真正需要登录权限的用户的角色。从 Postgres 8 开始，我们可以用 `create user` 来明确指出创建的角色是一个普通的用户账户，其默认赋予可登录权限；而 `create role` 所创建的角色被 Postgres 认为可能作为组角色，默认是不具登陆权限 (*cannot login*) 的^[6]，这也是 Postgres 中“用户”与“用户组”最核心的区别。

上图 (b) 和 (d) 示 PostgreSQL 13.3 中依然有“用户”的概念，但其只是作为可登录角色的别名^[7,8]。

```
CREATE USER [usr] = CREATE ROLE [usr] WITH LOGIN
```

尽管 Postgres 中的“组角色”也可以作为另一个组角色的成员角色，这种嵌套关系是无限层数的，但一般为了安全性，DBA 需要避免过多层组角色的继承以避免无意的为不该拥有某权限的角色授权。

5.2.2 权限授予

当 DBMS 收到一条请求后，在首先检查其语法无误后，DBMS 会检查发送请求的角色 (*role*) 及其是否有相应权限。如下图所示，如果发送 SQL 请求的角色没有对应权限，DBMS 将拒绝该请求。为角色授权的语法如下^[9]：

```
1 GRANT { { SELECT | INSERT | UPDATE | DELETE | TRUNCATE | REFERENCES | TRIGGER }
2       [, ...] | ALL [ PRIVILEGES ] }
3 ON { [ TABLE ] table_name [, ...]
4      | ALL TABLES IN SCHEMA schema_name [, ...] }
5 TO role_specification [, ...] [ WITH GRANT OPTION ]
```

```
postgres=# CREATE ROLE demogroup INHERIT; CREATE ROLE user1; CREATE USER user2;
CREATE ROLE
CREATE ROLE
postgres=# \du
```

Role name	List of roles Attributes	Member of
chris	Superuser, Create DB	{}
demogroup	Cannot login	{}
postgres	Superuser, Create role, Create DB, Replication, Bypass RLS	{}
user1	Cannot login	{}
user2		{}

```
postgres=# ALTER ROLE user1 WITH LOGIN;
ALTER ROLE
postgres=# \du
```

Role name	List of roles Attributes	Member of
chris	Superuser, Create DB	{}
demogroup	Cannot login	{}
postgres	Superuser, Create role, Create DB, Replication, Bypass RLS	{}
user1		{}
user2		{}

(a) 创建两个可登录角色、一个不可登录角色（组角色；*INHERIT* 关键词可省略）

rolname	rolsuper	rolinherit	rolcreaterole	rolcreatedb	rolcanlogin	rolreplication
1 chris	• true	• true	false	• true	• true	false
2 demogroup	false	• true	false	false	false	false
3 pg_execute_server_program	false	• true	false	false	false	false
4 pg_monitor	false	• true	false	false	false	false
5 pg_read_all_settings	false	• true	false	false	false	false
6 pg_read_all_stats	false	• true	false	false	false	false
7 pg_read_server_files	false	• true	false	false	false	false
8 pg_signal_backend	false	• true	false	false	false	false
9 pg_stat_scan_tables	false	• true	false	false	false	false
10 pg_write_server_files	false	• true	false	false	false	false
11 postgres	• true	• true	• true	• true	• true	• true
12 user1	false	• true	false	false	• true	false
13 user2	false	• true	false	false	• true	false

(b) `SELECT * FROM pg_roles;`

rolname
1 postgres
2 chris
3 user1
4 user2

username	usesysid
1 postgres	10
2 chris	16384
3 user1	24611
4 user2	25351

(c) `SELECT * FROM pg_roles WHERE rolcanlogin;`(d) `SELECT * FROM pg_user;`

```
[postgres@sngo]~% psql -h 127.0.0.1 -p 5432 -d postgres -U user1
WARNING: password file "/home/postgres/.pgpass" has group or world access
Password for user user1:
psql (13.3)
Type "help" for help.

postgres=> CREATE DATABASE demo;
ERROR: permission denied to create database
```

```
postgres=# ALTER ROLE user1 CREATEDB;
ALTER ROLE
postgres=# \du
```

Role name	List of roles Attributes	Member of
chris	Superuser, Create DB	{}
demogroup	Cannot login	{}
postgres	Superuser, Create role, Create DB, Replication, Bypass RLS	{}
user1	Create DB, Cannot login	{}
user2		{}

```
[postgres@sngo]~% psql -U user1 -h 127.0.0.1 -p 5432 -d postgres
WARNING: password file "/home/postgres/.pgpass" has group or world access
Password for user user1:
psql (13.3)
Type "help" for help.

postgres=> CREATE DATABASE demo1;
CREATE DATABASE
```

特别地，一个“对象”（schema、table 等）的所有者自然拥有该对象的所有权限。然而，作为对象的所有者并不意味着其是该对象的所有子对象的所有者（见下页图 (c)）。比如，如果一个角色创建并拥有一个数据库，而另一个角色在此数据库中创建了一个 schema，则该对象默认并没有权限去访问此 schema，但可以直接删除它。

只有权限的拥有者（其必须同时拥有该权限的 `grant` 权限）可以将权限授予其他角色，而有些权限只能由对象的拥有者持有，如 `DROP` 和 `ALTER`。另外，当授予一个角色时，可以通过添加 `WITH GRANT OPTION` 来允许该角色授予其他角色。

当创建一个新的数据库时，DBMS 会默认创建一个名为 *public* 的 schema，并将该 schema 的访问权授

予一个名为 `public` 的角色。所有新的用户和角色都被默认授予 `public` 角色中的所有权限，因此所有人都可以在此 `public schema` 中创建对象。

```
[postgres@sngo]~% psql -U user1 -d postgres -h 127.0.0.1
WARNING: password file "/home/postgres/.pgpass" has group
Password for user user1:
psql (13.3)
Type "help" for help.

postgres=> SELECT * FROM information_schema.schemata;
 catalog_name | schema_name | schema_owner | default_privileges
-----+-----+-----+-----
 postgres    | information_schema | postgres    |
 postgres    | public        | postgres    |
 postgres    | pg_catalog    | postgres    |
(3 rows)

postgres=> SET search_path TO public;
SET
postgres=> CREATE TABLE u1(id int);
CREATE TABLE
```

```
demo1=> SELECT * FROM information_schema.schemata;
catalog_name | schema_name | schema_owner | default_character_set_catalog | default_character_set_schema | default_character_set_name | sql_path
-----
demo1        | information_schema | postgres      |                               |                               |                               |
demo1        | public        | postgres      |                               |                               |                               |
demo1        | pg_catalog    | postgres      |                               |                               |                               |
(3 rows)

demo1=> CREATE SCHEMA u1; SET search_path TO u1; CREATE TABLE t1(name text); INSERT INTO t1 VALUES ('Chris'), ('Andrew');
CREATE SCHEMA
SET
CREATE TABLE
INSERT 0 2
demo1=> SELECT * FROM t1;
name
-----
Chris
Andrew
(2 rows)
```

(a) 对象的创建者自动拥有该对象的所有权限

```
demo1=> GRANT ALL PRIVILEGES ON DATABASE demo1 TO user2;
GRANT
demo1=> \q
[postgres@sngo]~% psql -U user2 -d demo1 -h 127.0.0.1
WARNING: password file "/home/postgres/.pgpass" has group or world access; permissions should be u=rw (0600) or less
Password for user user2:
psql (13.3)
Type "help" for help.

demo1=> SELECT * FROM u1.t1;
ERROR: permission denied for schema u1
LINE 1: SELECT * FROM u1.t1;
              ^
demo1=> CREATE SCHEMA u2; SET search_path TO u2; CREATE TABLE t2(id int); INSERT INTO t2 VALUES (1), (2); SELECT * FROM u2.t2;
CREATE SCHEMA
SET
CREATE TABLE
INSERT 0 2
demo1=> \q
[postgres@sngo]~% psql -U user2 -d demo1 -h 127.0.0.1
WARNING: password file "/home/postgres/.pgpass" has group or world access; permissions should be u=rw (0600) or less
Password for user user2:
psql (13.3)
Type "help" for help.

demo1=> SELECT * FROM information_schema.schemata;
catalog_name | schema_name | schema_owner | default_character_set_catalog | default_character_set_schema | default_character_set_name | sql_path
-----
demo1        | information_schema | postgres      |                               |                               |                               |
demo1        | public        | postgres      |                               |                               |                               |
demo1        | pg_catalog    | postgres      |                               |                               |                               |
demo1        | u1            | user1         |                               |                               |                               |
(4 rows)

demo1=> \q
[postgres@sngo]~% psql -U user2 -d demo1 -h 127.0.0.1
WARNING: password file "/home/postgres/.pgpass" has group or world access; permissions should be u=rw (0600) or less
Password for user user2:
psql (13.3)
Type "help" for help.

demo1=> SELECT * FROM information_schema.schemata;
catalog_name | schema_name | schema_owner | default_character_set_catalog | default_character_set_schema | default_character_set_name | sql_path
-----
demo1        | information_schema | postgres      |                               |                               |                               |
demo1        | public        | postgres      |                               |                               |                               |
demo1        | pg_catalog    | postgres      |                               |                               |                               |
demo1        | u2            | user2         |                               |                               |                               |
(4 rows)
```

(b) 角色被授予某数据库的所有权限后，依然不能访问特定的 schema（除 public 外，由其他角色创建），但可以对数据库进行其他操作

```
[postgres@sngo]~% psql -U user1 -d demo1 -h 127.0.0.1
WARNING: password file "/home/postgres/.pgpass" has group or world access; permissions should be u=rw (0600) or less
Password for user user1:
psql (13.3)
Type "help" for help.

demo1=> SELECT * FROM information_schema.schemata;
catalog_name | schema_name | schema_owner | default_character_set_catalog | default_character_set_schema | default_character_set_name | sql_path
-----
demo1        | information_schema | postgres      |                               |                               |                               |
demo1        | public        | postgres      |                               |                               |                               |
demo1        | pg_catalog    | postgres      |                               |                               |                               |
demo1        | u1            | user1         |                               |                               |                               |
(4 rows)

demo1=> \q
[postgres@sngo]~% psql -U user2 -d demo1 -h 127.0.0.1
WARNING: password file "/home/postgres/.pgpass" has group or world access; permissions should be u=rw (0600) or less
Password for user user2:
psql (13.3)
Type "help" for help.

demo1=> SELECT * FROM information_schema.schemata;
catalog_name | schema_name | schema_owner | default_character_set_catalog | default_character_set_schema | default_character_set_name | sql_path
-----
demo1        | information_schema | postgres      |                               |                               |                               |
demo1        | public        | postgres      |                               |                               |                               |
demo1        | pg_catalog    | postgres      |                               |                               |                               |
demo1        | u2            | user2         |                               |                               |                               |
(4 rows)
```

(c) user2 在 user1 所拥有对数据库中创建的 schema 对 user1 不可见

```
[postgres@sngo]~% psql -U user1 -d postgres -h 127.0.0.1
WARNING: password file "/home/postgres/.pgpass" has group or world access; permissions should be u=rw (0600) or less
Password for user user1:
psql (13.3)
Type "help" for help.

postgres=> DROP DATABASE demo1;
DROP DATABASE
```

(d) user1 未被授予 schema u2 的相关权限，但因其拥有整个数据库，因此可以直接将其 drop

5.2.3 撤回权限

收回权限的语法与授权的语法类似：

```

1 REVOKE [ GRANT OPTION FOR ]
2   { { SELECT | INSERT | UPDATE | DELETE | TRUNCATE | REFERENCES | TRIGGER }
3     [, ...] | ALL [ PRIVILEGES ] }
4   ON { [ TABLE ] table_name [, ...]
5         | ALL TABLES IN SCHEMA schema_name [, ...] }
6   FROM role_specification [, ...]
7   [ CASCADE | RESTRICT ]

```

```

postgres=# REVOKE ALL ON TABLE u1 FROM user1;
REVOKE

```

在 Postgres 中较为特殊的一点是，删除角色时，一般不能直接删除，而要先将该角色被授予的所有权限收回，并转移其拥有的对象的所有权。

```

postgres=# DROP ROLE user1;
ERROR:  role "user1" cannot be dropped because some objects depend on it
DETAIL:  owner of table u1
2 objects in database grp
postgres=# ALTER TABLE u1 OWNER TO user2;
ALTER TABLE
postgres=# DROP ROLE user1;
ERROR:  role "user1" cannot be dropped because some objects depend on it
DETAIL:  2 objects in database grp

```

5.2.4 组角色

在前两小节的演示角色的图 (a) 中我们创建了一个名为 *demogroup* 的组角色，现用角色 *postgres* 创建一个数据库，并将其所有权限赋予给该组角色。注意我们至始至终没有给 *user1* 授予权限。接下来，通过

GRANT [group_role] TO [login_role];

命令将 *user1* 和 *user2* 加入 *demogroup* 的成员角色。下图演示了进入组角色的 *user1* 自动继承了组中所有权限，我们得以在不直接对 *user1* 授权的情况下使 *user1* 有权操作该数据库。

```

postgres=# CREATE DATABASE grp; GRANT ALL PRIVILEGES ON DATABASE grp TO demogroup;
CREATE DATABASE
GRANT
postgres=# GRANT demogroup TO user1; GRANT demogroup TO user2;
GRANT ROLE
GRANT ROLE
postgres=# \du

```

Role name	List of roles Attributes	Member of
chris	Superuser, Create DB	{}
demogroup	Cannot login	{}
postgres	Superuser, Create role, Create DB, Replication, Bypass RLS	{}
user1	Create DB	{demogroup}
user2		{demogroup}

```
[postgres@sngo]~% psql -U user1 -d grp -h 127.0.0.1
WARNING: password file "/home/postgres/.pgpass" has group or world access; permission denied
Password for user user1:
psql (13.3)
Type "help" for help.

grp=> CREATE SCHEMA uis1; CREATE TABLE y(id int); INSERT INTO y VALUES (1);
CREATE SCHEMA
CREATE TABLE
INSERT 0 1
```

要删除一个组角色，执行 `DROP ROLE [group_role]` 命令即可。在删除该组角色之后，它与其成员角色之间的关系将被立即撤销，而成员角色本身不受影响。需要注意的是，在删除前任何属于该组角色的对象都必须先被删除或者将对象的所有者转移给其它角色，与此同时，任何赋予该组角色的权限也都必须被撤销^[10]。

5.2.5 文件系统用户权限管理

文件系统可以通过改变文件的权限（以 UNIX 系 OS 为例，一般我们关注读、写、可执行权限）来达到类似的效果，通过系统自带的用户机制可以为不同用户、针对某个文件，分配不同的权限。

```
(base) chris /home/chris ll
total 12K
drwx----- 22 chris chris 4.0K Nov  8 21:49 anaconda3
-rw-----  1 root  root  5.3K Nov 16 11:48 course_info.json
(base) chris /home/chris cat course_info.json
cat: course_info.json: Permission denied
(base) chris /home/chris ll
total 12K
drwx----- 22 chris chris 4.0K Nov  8 21:49 anaconda3
-rw-rw-rw-  1 root  root  5.3K Nov 16 11:48 course_info.json
(base) chris /home/chris cat course_info.json
[
  {
    "totalCapacity": 28,
    "courseId": "GE232",
    "prerequisite": null,
```

上图演示在 CentOS 中使用 root 账号创建文件，并执行 `chmod 600` 命令使文件仅创建者可读写（如第一个 `ll` 所示），此时登录另一账号尝试读取文件，权限不足被拒；再次执行 `chmod 666` 对所有用户开放读写权限（如第二个 `ll` 所示），再次读取文件则成功显示内容。

5.3 数据库索引及文件 IO

5.3.1 数据库索引

对于没有 index 的表的查询，DBMS 需要将整行数据取出并在其中少数几个字段应用筛选器，这种操作是低效的。而为了能直接定位所需的记录，DBMS 设计了索引这种与文件相关联的附加的结构。常用的索引形式为 B⁺ 树及散列索引。

当查询的限定字段被某一索引所包括时，DBMS 首先查找索引，找出相应记录所在的磁盘块，再取出相应的（存储数据量较小的）磁盘块并在其中高效搜索^[4]。但索引并非越多越好：一方面，创建索引和动态维护索引要耗费时间；另一方面，索引需要占物理空间，除了数据表占数据空间之外，每一个索引还要占一定的物理空间，如果要建立聚簇索引，那么需要的空间就会更大。

在以下情况下使用索引能带来较优的效果^[4]:

- ① 频繁需要搜索的列;
- ② 主键列、唯一性约束、外键列等, 强制该列的唯一性和组织表中数据的排列结构 (PostgreSQL 会自动为这些列设置索引);
- ③ 连接表或检查被用作连接 (join) 的列, 即一般指外键列;
- ④ 经常需要进行范围搜索的列, 可利用索引已排序的特性对数值进行连续搜索;
- ⑤ 经常需要排序的列同样可利用索引已排序的特性减少排序步骤的耗时;
- ⑥ 频繁作为限定条件的字段组合 (列)。

而以下情况使用索引的效果并不佳:

- ① 查询中很少使用或者参考的列不应该创建索引 (对查询速度的优化场景较少, 而 DBMS 多了储存与维护索引的时间、空间负担);
- ② 数据量较少的列不需要设置索引;
- ③ 连接表或检查被用作连接 (join) 的列, 即一般指外键列;
- ④ Text, image, blob 等数据类型的列不应该增加索引, 因为这些列的数据量较大或取值很少;
- ⑤ 需要频繁修改而较少需要检索的表: 修改性能和检索性能互相矛盾, 当增加索引时, 会提高检索性能, 但是会降低修改性能, 当减少索引时, 会提高修改性能, 降低检索性能。

```
QUERY PLAN
1  Index Scan using student_pkey on student  (cost=0.29..8.31 rows=1 width=22)
2    Index Cond: (id = 11080020)
```

SELECT * FROM project1.student WHERE id = 11080020; 这里 id 作为主键列已经设有 index

以下查询所有与海洋相关的院系所开设的课程及教师。此时 join 中用到的连接列均有数据库默认加上 index, 键入 EXPLAIN ANALYSE 命令分析可见大量 *hash join*, *hash cond*, *index scan*, *index cond* 等, 查询耗时 0.368ms。

```
1  SELECT c.name, t.name
2  FROM project1.course c
3       JOIN project1.department d ON c.dept = d.id
4       JOIN project1.class cls ON cls.course = c.cid
5       JOIN project1.class_teacher ct ON ct.class = cls.id
6       JOIN project1.teacher t ON t.id = ct.teacher
7  WHERE d.name ~ '海洋'
8  ORDER BY c.name, t.name;
```

```

QUERY PLAN
1 Sort (cost=31.32..31.33 rows=3 width=29) (actual time=0.307..0.311 rows=11 loops=1)
2   Sort Key: c.name, t.name
3   Sort Method: quicksort Memory: 25kB
4   → Nested Loop (cost=21.96..31.30 rows=3 width=29) (actual time=0.100..0.272 rows=11 loops=1)
5     → Nested Loop (cost=21.81..30.73 rows=3 width=24) (actual time=0.087..0.227 rows=11 loops=1)
6       → Nested Loop (cost=21.54..29.99 rows=2 width=24) (actual time=0.088..0.192 rows=12 loops=1)
7         → Hash Join (cost=21.26..29.53 rows=1 width=26) (actual time=0.068..0.137 rows=12 loops=1)
8           Hash Cond: (c.dept = d.id)
9           → Seq Scan on course c (cost=0.00..7.38 rows=338 width=30) (actual time=0.014..0.045 rows=338 loops=1)
10            → Hash (cost=21.25..21.25 rows=1 width=4) (actual time=0.041..0.042 rows=1 loops=1)
11              Buckets: 1024 Batches: 1 Memory Usage: 9kB
12              → Seq Scan on department d (cost=0.00..21.25 rows=1 width=4) (actual time=0.032..0.038 rows=1 loops=1)
13                Filter: ((name)::text ~ '海洋'::text)
14                Rows Removed by Filter: 22
15            → Index Scan using class_unk on class cls (cost=0.28..0.44 rows=2 width=10) (actual time=0.004..0.004 rows=1 loops=1)
16              Index Cond: ((course)::text = (c.cid)::text)
17          → Index Only Scan using cls_tec_unk on class_teacher ct (cost=0.28..0.36 rows=1 width=8) (actual time=0.002..0.002 rows=1 loops=1)
18            Index Cond: (class = cls.id)
19            Heap Fetches: 11
20        → Index Scan using teacher_pkey on teacher t (cost=0.15..0.19 rows=1 width=13) (actual time=0.002..0.002 rows=1 loops=1)
21          Index Cond: (id = ct.teacher)
22 Planning Time: 0.886 ms
23 Execution Time: 0.368 ms

```

下面通过对比能体现索引带来了极高的查询效率提升。对于查询语句 **SELECT class FROM schedule WHERE weekday = 5**，原始的无索引查询（如左下图）进行了序列扫描，需要取出表中所有信息，耗时 0.281s，而通过 **CREATE INDEX ON project1.schedule (weekday)** 为 WHERE 限定列加上索引后，查询耗时仅为 0.147s（如右下图），相比无索引查询节省了 47.6% 的时间。

```

QUERY PLAN
1 Seq Scan on schedule (cost=0.00..29.11 rows=142 width=10)
2   Filter: (weekday = 5)
3   Rows Removed by Filter: 827
4 Planning Time: 12.815 ms
5 Execution Time: 0.281 ms

```

```

QUERY PLAN
1 Bitmap Heap Scan on schedule (cost=5.25..24.83 rows=142 width=4) (actual time=0.068..0.106 rows=142 loops=1)
2   Recheck Cond: (weekday = 5)
3   Heap Blocks: exact=17
4   → Bitmap Index Scan on schedule_weekday_idx (cost=0.00..5.21 rows=142 width=0) (actual time=0.049..0.049 rows=142 loops=1)
5     Index Cond: (weekday = 5)
6 Planning Time: 0.269 ms
7 Execution Time: 0.147 ms

```

5.3.2 文件索引（数据结构）

使用 Python 对无索引的课程信息文件十次检索学分大于 2 的课程，耗时 0.525s：

```

1 def fs_query() -> set:
2     result_set = set()
3     for _ in range(10):
4         result_set.clear()
5         for index, c in cif.iterrows():
6             if c['courseCredit'] > 2:
7                 result_set.add(c['courseId'])
8     return result_set

```

下面为相应条目创建索引³ [11]

```

1 Index idx = new Index();
2 idx.index("courseCredit", "course_info.json", "idx_course_info.bpt");
3 idx.searchindex("idx_course_info.bpt", "courseCredit > 2", " ");

```

调用该封装好的函数如同使用 SQL 一样方便，首先需要对对应的字段创建索引，然后解析此索引并进行搜索，并加工 searchindex 函数返回的原始结果（相对于 SELECT *），十次搜索耗时 0.319s，但考虑到 JVM 和 Python interpreter 的性能差异，该结果仅供参考。

³使用了 GitHub 上的 B+ 树代码实现：adigan1310 / BPlus-Tree-Indexing

5.4 跨平台性能对比

File system ① CentOS 7.9.2009 (Python 3.8.8) ② macOS 12.0.1 (Python 3.8.10)

RDBMS ① PostgreSQL 13.3 ② MySQL 5.6.50 ③ SQLite 3.34.0

数据均以相同结构存储在数据库中，而文件系统部分考虑到 Java 历史久远，JVM 在各平台下的表现高度均一，仅能体现不同评测环境的硬件区别，故使用 Python 进行对比，使用 Pandas 库读取 csv 文件。以下是查询所有修完 CS205 的学生信息的平均耗时（文件系统重复十次查询取平均，数据库不预热；下表中数据库项均为 execution + fetching，由于使用的均为云数据库，fetching time 受网络环境延迟影响数值较大）：

	macOS (Python)	CentOS (Python)	PostgreSQL	MySQL	SQLite
平均耗时 (ms)	1139	1271	345 + 112	317 + 107	359 + 120

注意到 PostgreSQL 在查询速度上略低于 MySQL，这是因为 PostgreSQL 和 MySQL 均是基于 MVCC 机制 (Multi-Version Concurrency Control) 来保证事务的原子性和隔离性的，该模型将争用降低到最低限度以期减少操作的阻塞。PostgreSQL 的实现和 MySQL 基于的 InnoDB 引擎的实现是目前公认的两大实现 MVCC 的方法。其主要区别在于 Postgres 在查询开始时为了避免脏读而存储的旧版本，其只是在元组的较旧版本上更新 xmax，因此较旧版本的大小与相应的插入记录相同。而对于 InnoDB，存储在 Undo 段中的对象版本通常小于相应的插入记录。这是因为仅将更改的值（即差分）写入 UNDO 日志。^[12]而轻量级的 SQLite 相较于另外两种数据库，更接近于文件系统。文件系统的查询速度有细微差异，这可能是由于 Python interpreter 在不同系统上的实现略微不同以及不同系统对文件存储策略的不同造成的。

本次测试中也注意到，数据库在第一次查询时耗时（345ms）远大于后续几次进行相同查询（30ms），这是由于 DBMS 的缓存策略相比起文件系统中实现的更加健全高效，在第一次将数据搬入缓存或内存中后，后续查询能有较高的缓存命中率。

References

- [1] "VARCHAR as foreign key/primary key in database good or bad?" 01 2010. [Online]. Available: <https://stackoverflow.com/questions/2103322/varchar-as-foreign-key-primary-key-in-database-good-or-bad>
- [2] J. M, "一条 SQL 的执行过程详解," 12 2020. [Online]. Available: <https://www.cnblogs.com/mengxin/p/14045520.html>
- [3] T. Dong, "5 分钟理解数据库死锁," 05 2021. [Online]. Available: <https://blog.csdn.net/horses/article/details/116503824>
- [4] A. Silberschatz, H. Korth, and S. Sudarshan, *Database System Concepts*, 6th ed. McGraw-Hill Education, 2010.
- [5] Wikipedia contributors, "ACID," 10 2021. [Online]. Available: <https://en.wikipedia.org/wiki/ACID>
- [6] R. Obe and L. Hsu, *PostgreSQL: Up and Running: A Practical Guide to the Advanced Open Source Database*, 3rd ed. O'Reilly Media, 2017.
- [7] "Database Roles and Privileges," 01 2012. [Online]. Available: <https://www.postgresql.org/docs/8.1/user-manag.html>
- [8] "What is the difference between a user and a role?" 12 2014. [Online]. Available: <https://stackoverflow.com/questions/27709456/what-is-the-difference-between-a-user-and-a-role>
- [9] "Privileges," 08 2021. [Online]. Available: <https://www.postgresql.org/docs/13/ddl-priv.html>
- [10] M. Turing, "postgresql-revoke 回收权限及删除角色," 02 2021. [Online]. Available: <https://www.cnblogs.com/zhangfx01/p/14367594.html>

- [11] Adigan1310, "BPlus-Tree-Indexing: Implementation of B+ Tree Indexing of a file using memory blocks," 2017. [Online]. Available: <https://github.com/adigan1310/BPlus-Tree-Indexing>
- [12] S. Blue, "PostgreSQL 和 Mysql 的 MVCC 实现机制的差异对比," 02 2020. [Online]. Available: <https://zhuanlan.zhihu.com/p/107993134>

Disclaimer: All references used in this project are only for the purpose of guiding the direction of exploration and ensuring correctness, any code appearing in the referenced web pages has not been used directly, unless explicitly stated.