# PyShare 03

👤 lyj8512@126.com

# Python ≠ 爬虫

# Python 语言,工具

Python ●

✅ **科学计算**

✅ **机器学习,深度学习,AI领域**

✅ **数据清洗,处理,分析,统计**

😄 爬虫

💻 Scalable Web Crawler ::**通用爬虫**

Focused Crawler & Topical Crawler ::**聚焦爬虫**

Incremental Web Crawler ::**增量式爬虫**

Deep Web Crawler ::**深层网页爬虫**

🎯 ！期望能获取到有研究价值的数据

♥ *图像,自然语言处理*

♥ *其它,你能想到的....*

安装第三方库

Pip install pyquery

Pip install requests

Pip install aiothhp

# 代码实现,用什么?

*编写,调试,运行…*

# 下载安装

# Pycharm

- PyCharm 介绍
- 项目实现
  - 通过PyCharm,创建python 爬虫项目
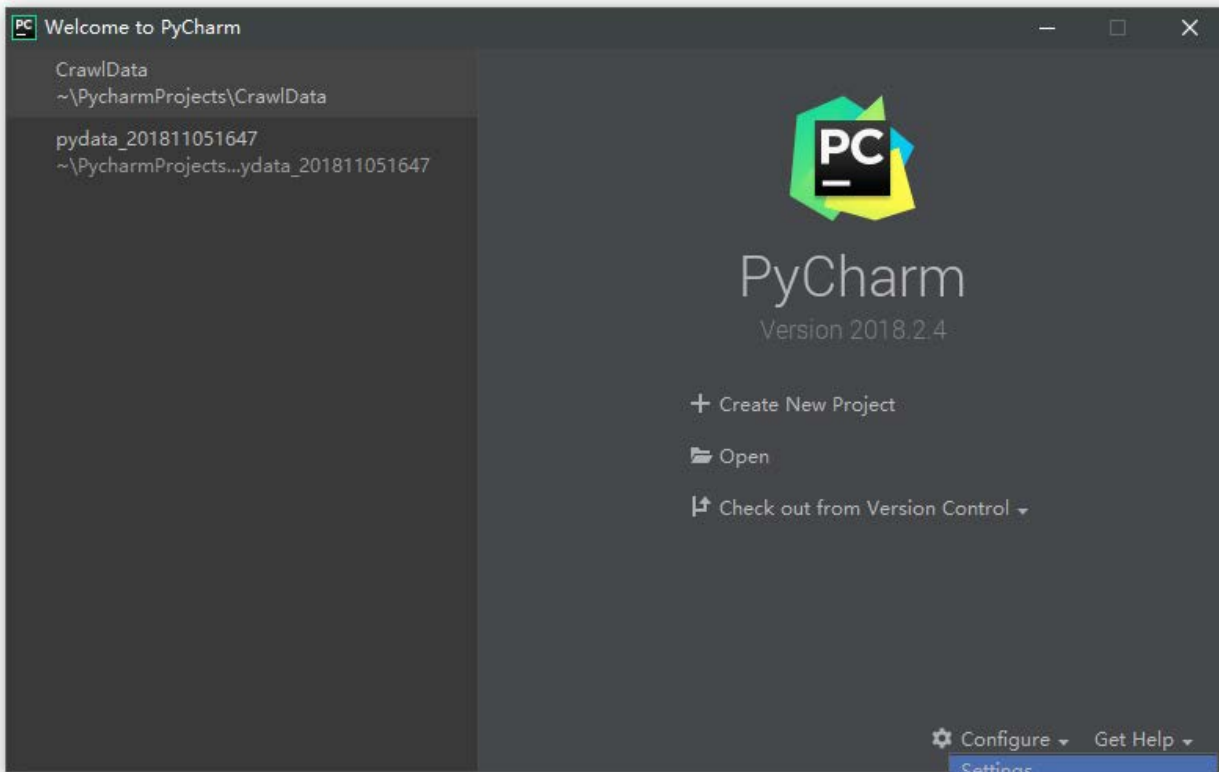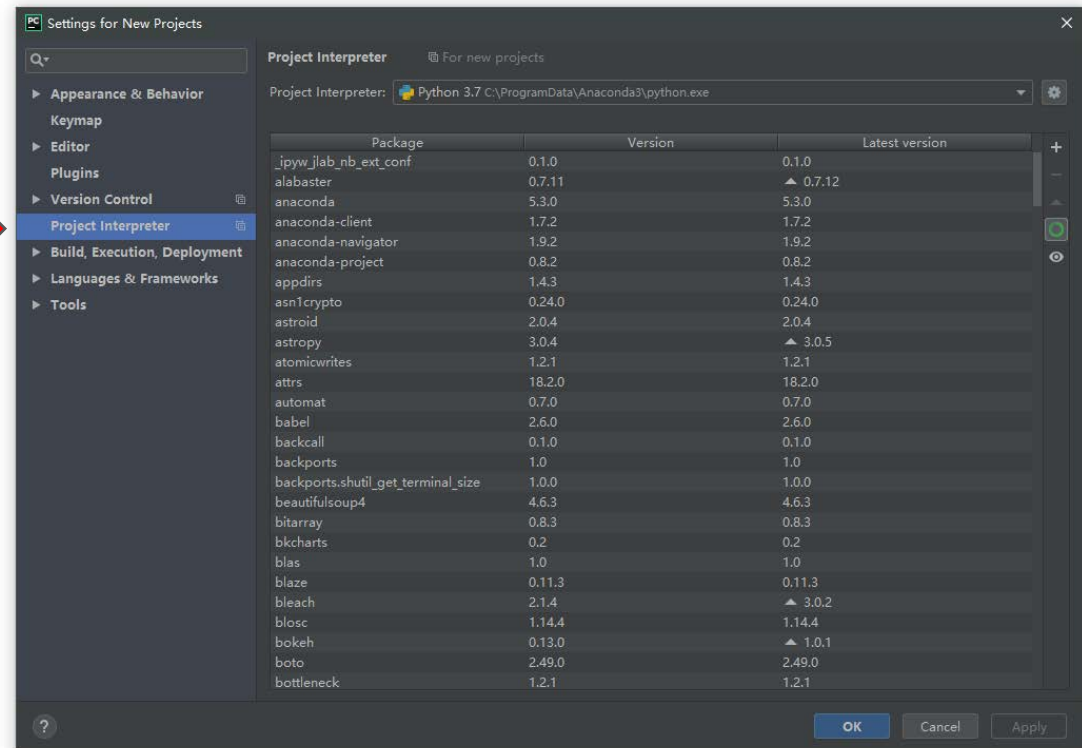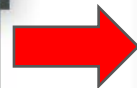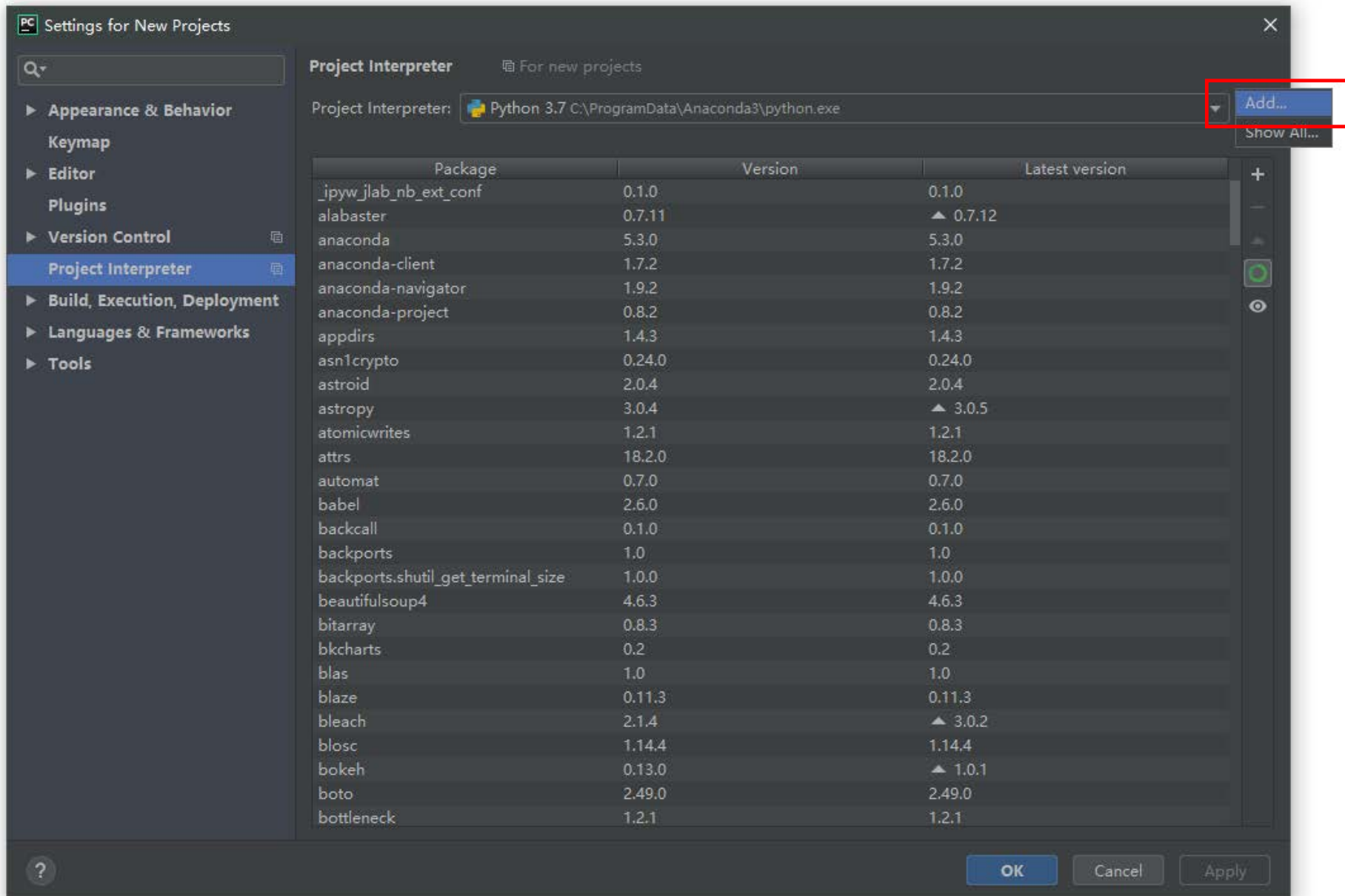  - 配置anaconda解释环境
- PyCharm常用功能
  - 理解调试的作用
  - 快捷键
    - 格式化代码
    - 代码定位,查找,重构等
    - 运行,调试

# Pycharm 配置

*配置加载需要时间,…*

**Settings for New Projects**

Project Interpreter    📋 For new projects

Project Interpreter:    🐍 Python 3.7 C:\ProgramData\Anaconda3\python.exe    ▼    Add...

Show All...

- ▶ Appearance & Behavior
- Keymap
- ▶ Editor
- Plugins
- ▶ Version Control
- **Project Interpreter**
- ▶ Build, Execution, Deployment
- ▶ Languages & Frameworks
- ▶ Tools

| Package | Version | Latest version |
| --- | --- | --- |
| _ipyw_jlab_nb_ext_conf | 0.1.0 | 0.1.0 |
| alabaster | 0.7.11 | ▲ 0.7.12 |
| anaconda | 5.3.0 | 5.3.0 |
| anaconda-client | 1.7.2 | 1.7.2 |
| anaconda-navigator | 1.9.2 | 1.9.2 |
| anaconda-project | 0.8.2 | 0.8.2 |
| appdirs | 1.4.3 | 1.4.3 |
| asn1crypto | 0.24.0 | 0.24.0 |
| astroid | 2.0.4 | 2.0.4 |
| astropy | 3.0.4 | ▲ 3.0.5 |
| atomicwrites | 1.2.1 | 1.2.1 |
| attrs | 18.2.0 | 18.2.0 |
| automat | 0.7.0 | 0.7.0 |
| babel | 2.6.0 | 2.6.0 |
| backcall | 0.1.0 | 0.1.0 |
| backports | 1.0 | 1.0 |
| backports.shutil_get_terminal_size | 1.0.0 | 1.0.0 |
| beautifulsoup4 | 4.6.3 | 4.6.3 |
| bitarray | 0.8.3 | 0.8.3 |
| bkcharts | 0.2 | 0.2 |
| blas | 1.0 | 1.0 |
| blaze | 0.11.3 | 0.11.3 |
| bleach | 2.1.4 | ▲ 3.0.2 |
| blosc | 1.14.4 | 1.14.4 |
| bokeh | 0.13.0 | ▲ 1.0.1 |
| boto | 2.49.0 | 2.49.0 |
| bottleneck | 1.2.1 | 1.2.1 |

? 

OK    Cancel    Apply

untitled [C:\Users\ii\PycharmProjects\untitled] - ...\test.py [untitled] - PyCharm

File  Edit  View  Navigate  Code  Refactor  Run  Tools  VCS  Window  Help

untitled  test.py

Project

untitled  C:\Users\ii\PycharmProjects\untitled
    test.py
External Libraries
Scratches and Consoles

1    im

import
__import__(name, globals, locals, fr…   builtins
Press Ctrl+. to choose the selected (or first) suggestion and insert a dot afterwards >>

为什么要用
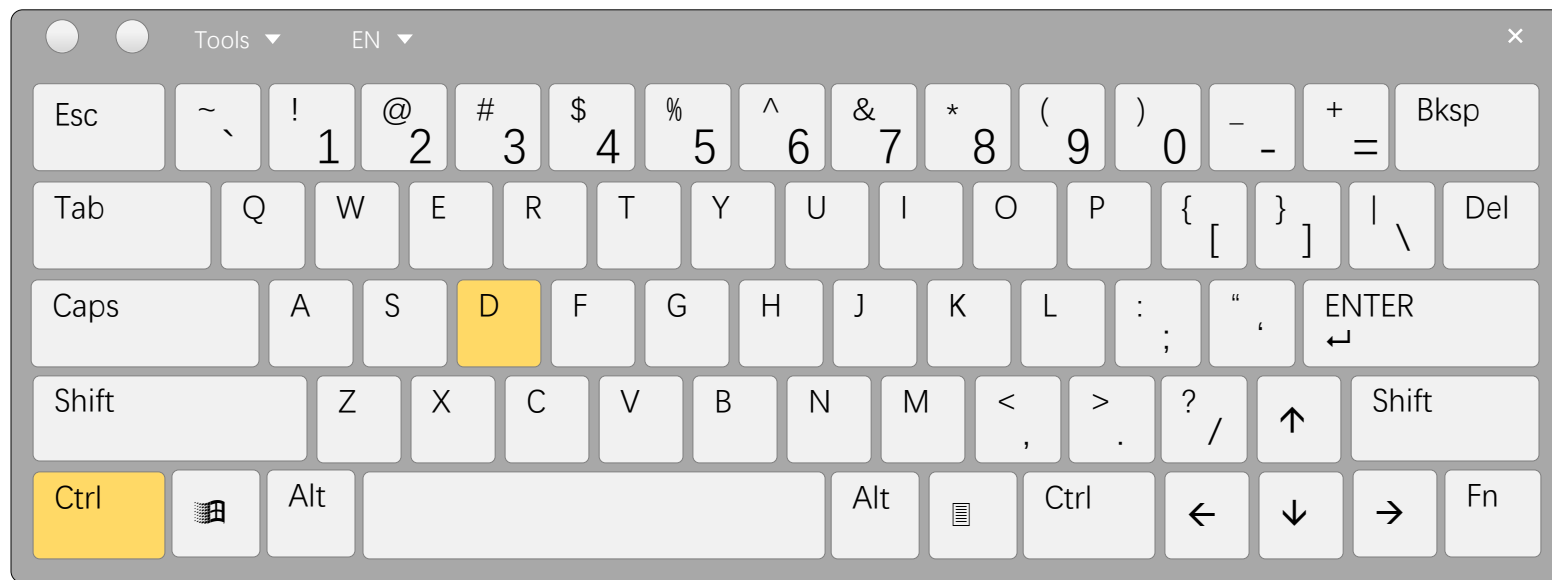PyCharm

提供(自检,提醒,智能补全)功能…

专注代码逻辑实现,…

```
1  import keyword
2
3  for i in range(10):
4      print(i)
5
6
```
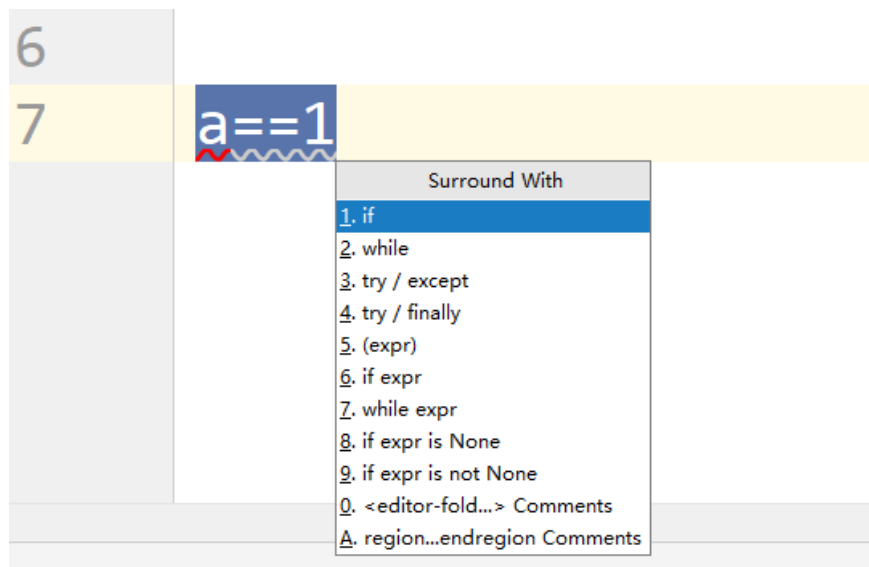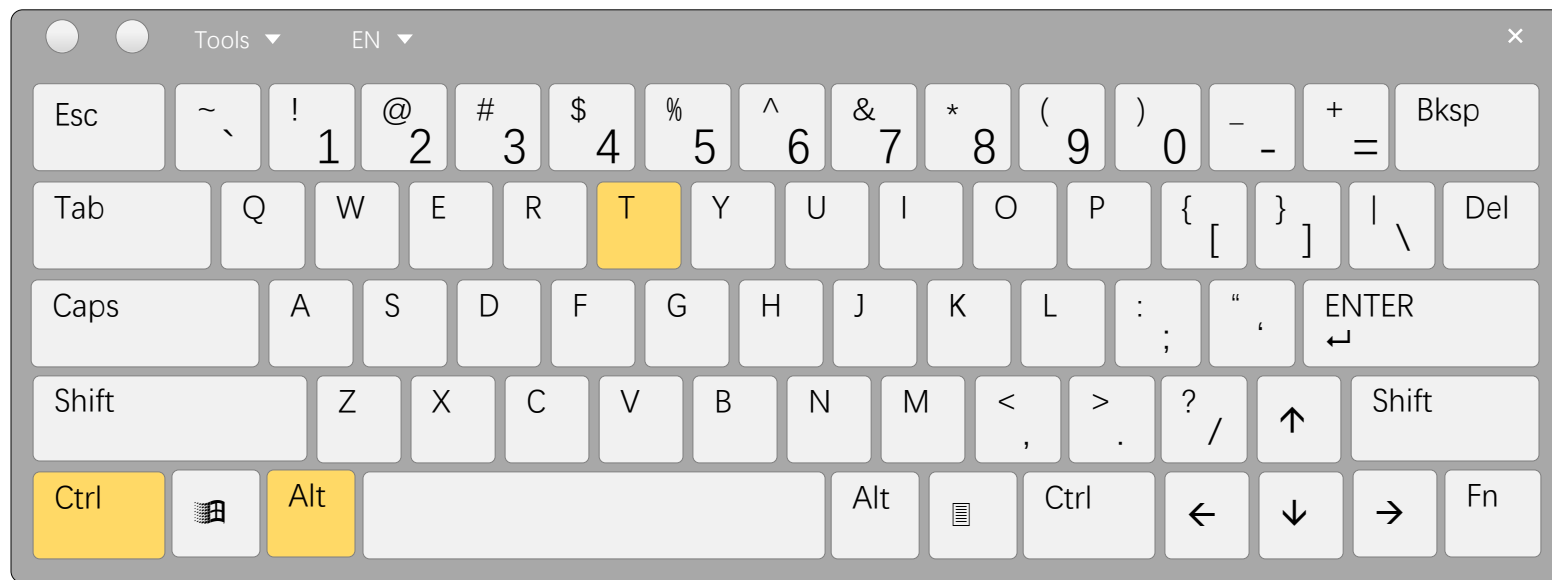
光标任意位置,切换一行,…

```
1  import keyword
2
3  for i in range(10):
4      print(i)
5      print(i)
```
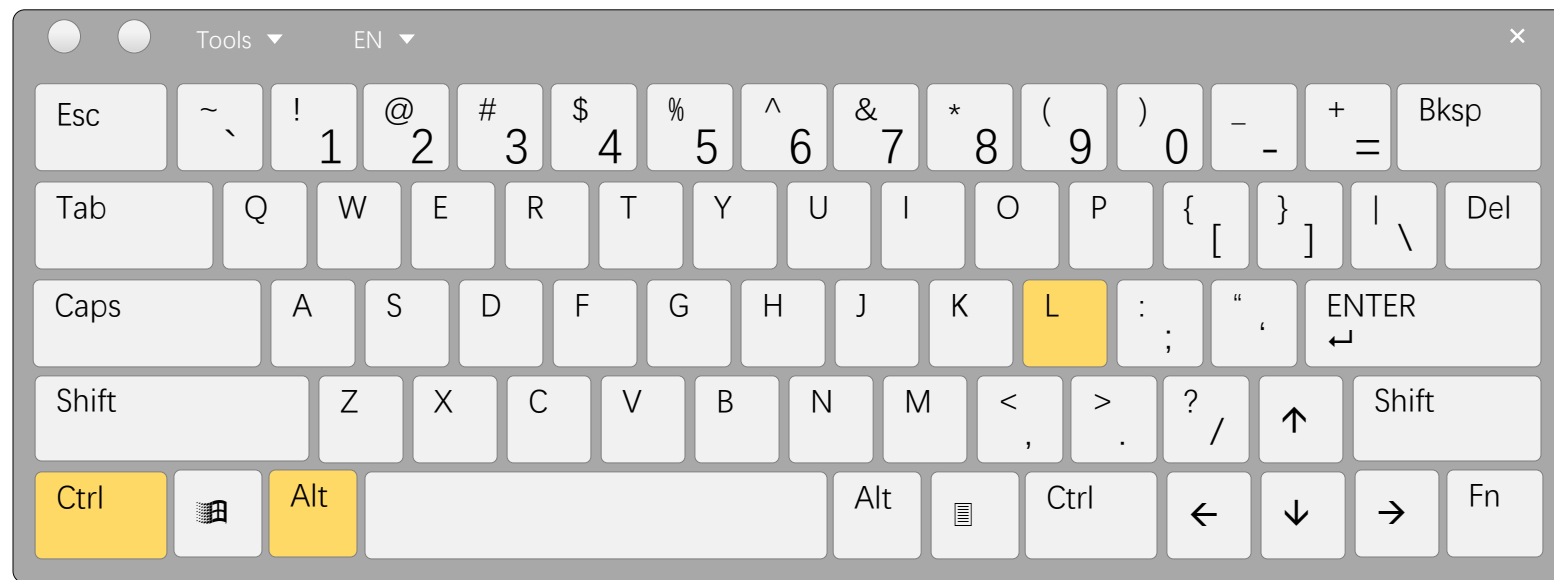
复制当前行,…

对当前行进行流程控制语句补全…

```
6
7   a=1
8   b=1
9   c=1
10  d=1
11
```

有波浪提示的时候,表示代码编写规范,…
用快捷键格式化格式即可,

```
7   # a = 1
8   # b = 1
9   # c = 1
10  # d = 1
11
```

代码注释切换,

更多快捷键:

https://blog.csdn.net/pipisorry/article/details/39909057

# PyCharm 个性化设置

设置字体大小
CTRL+鼠标滚轮控制

| Settings | | ✕ |
|---|---|---|

Editor › Color Scheme

> Appearance & Behavior
  Keymap
∨ Editor
  > General
    Font
    Color Scheme
  > Code Style
    Inspections
    File and Code Templates
    File Encodings
    Live Templates
    File Types
  > Emmet
    Images
    Intentions
    Language Injections
    Spelling
    TextMate Bundles
    TODO
  Plugins
> Version Control
> Project: untitled
> Build, Execution, Deployment
> Languages & Frameworks

Scheme: Default

**Default**
**Darcula**
Github
**Monokai**
Twilight
WarmNeon

OK    Cancel    Apply

设置主题

# 代码实现回顾
## 往期内容

### python 基础
- 列表 list
  - 索引&切片
  - 增删改查
- for 循环语句
- 如何(定义 & 调用)函数,

*http://www.runoob.com/python3/python3-tutorial.html*

### 简单爬虫

#### requests 网页请求库
- *respone=get(url)*
  *respone=get(url,headers=header)*
- *respone.text*

*http://docs.python-requests.org/en/master/*

#### PyQuery 网页解析库
- 理解标签 & 标签的属性值
- 理解CSS选择器
- 通过CSS选择器查找PyQuery对象
- 解析PyQuery对象-(内容获取)
  - items()
  - text()

*https://pythonhosted.org/pyquery/*

# PyCharm

练习:

创建列表,索引,切片
创建函数,调用函数…

# PyCharm

练习:简单爬虫实现

任务:猫眼网页:正在热映电影榜单
http://maoyan.com/films

`<code>` PyShare D3-00.py

# 猫眼网页:

*http://maoyan.com/films*



```python
import requests
from pyquery import PyQuery as pq

url = "http://maoyan.com/films"

# requests 请求头
headers = {
    "User-Agent": "Mozilla/5.0 (Windows NT 10.0; Win64; x64) AppleWebKit/537.36 "
                  "(KHTML, like Gecko) Chrome/70.0.3538.67 Safari/537.36"

}


def gethtml(url):
    temp = requests.get(url, headers=headers)
    return temp.text


def parserhtml(html):
    doc = pq(html)
    a = doc('.channel-detail.movie-item-title').items()

    for i in a:
        print(i.text())


if __name__ == '__main__':
    a = gethtml(url)
    parserhtml(a)
```

*Requests :*

*http://docs.python-requests.org/en/master/*


*PyQuery :*

*https://pythonhosted.org/pyquery/*

线程是什么

执行A事时,有等待返回的操作时,处理器停止(按顺序执行)

数据爬取随网页数量线性递增

协程是什么

执行A事件时,有等待返回的操作时,先挂起A事件,执行B事件,当A事件有返回值时,执行A(异步执行)

异步执行,节省时间,速度快

单线程 & 协程

Thread #1
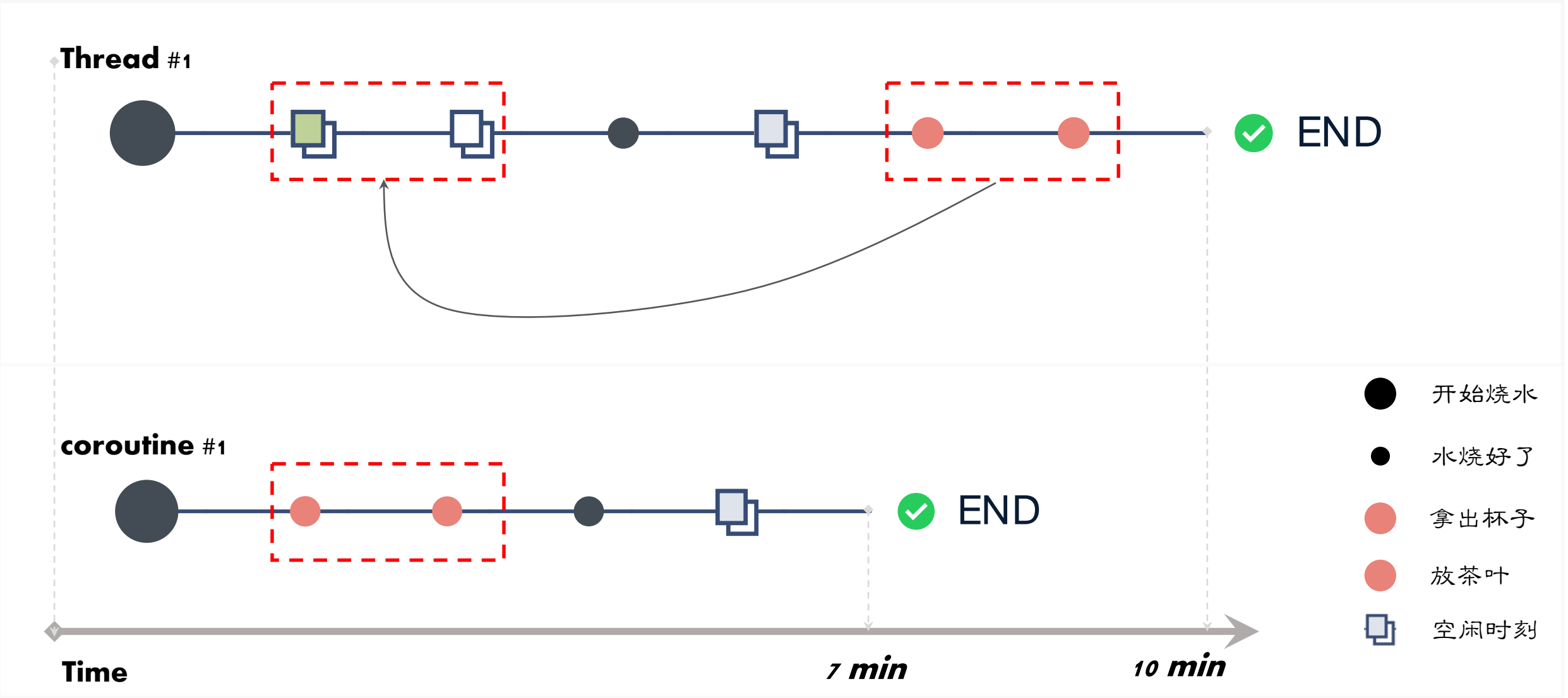
END

coroutine #1

END

Time

7 min

10 min

开始烧水

水烧好了

拿出杯子

放茶叶

空闲时刻

# 理解线程



Thread # 1

Time

- 开始烧水
- 水烧好了
- 可切换节点
- 空闲时刻

# 理解协程



**Thread** #1 ... **coroutine** #1 ... END

Time 7 min 10 min

- 开始烧水
- 水烧好了
- 拿出杯子
- 放茶叶
- 空闲时刻

# Coroutine
## 协程(异步爬虫)

### asycnio
### python 自带库

理解async 函数方法

创建协程 coroutine

> 直接将方法注册到事件中,
> loop.create_task()
> asyncio.ensure_future()

asyncio.get_event_loop()
创建协程的循环事件

loop.run_until_complete(asyncio.gather(*list))
loop.run_until_complete(asyncio.wait(list))
运行协程事件,获取结果

https://docs.python.org/3/library/asyncio.html

### asiohttp
### 第三方异步requests请求库

创建async 请求会话

await 有什么用:
异步爬虫中对有等待返回操作的对象进行挂起

与asycnio调用类似

https://aiohttp.readthedocs.io/en/stable/client_quickstart.html#make-a-request

*asycnio :*

*https://docs.python.org/3/library/asyncio.html*

*aiohttp :*

*https://aiohttp.readthedocs.io/en/stable/client_quickstart.html#make-a-request*

Process ⊗     *mulitProcess*
*多进程*

         *理解多线程*

*multThread*
*多线程*     *多线程解析html*
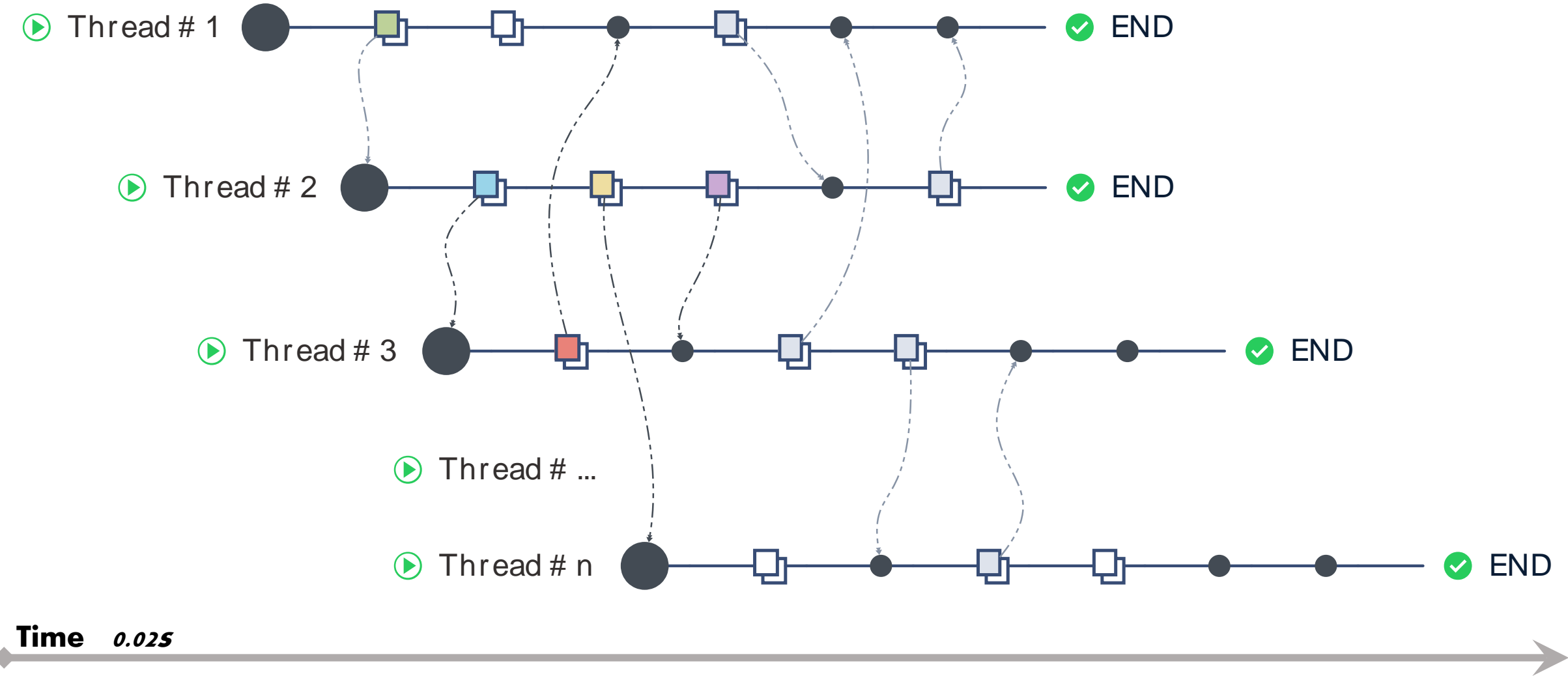
Thread ⊗

*线程锁*     *数据存取有关*

*线程池*     *线程执行数量有关*

**按进度调整**

**Concurrency** is about dealing with lots of things at once.

**Parallelism** is about doing lots of things at once.

# 理解多线程



Thread # 1      END

Thread # 2      END

Thread # 3      END

Thread # ...

Thread # n      END

**Time**   *0.02s*

# 理解进程



prosess # 1 ........ END

prosess # 2 ........ END

Time 0.02s