

译

4

Python Django性能测试与优化指南

2017年12月2日 15:43:35

Guide to Performance Testing and Optimization With Python and Django

LIAN GULEA

惊寒

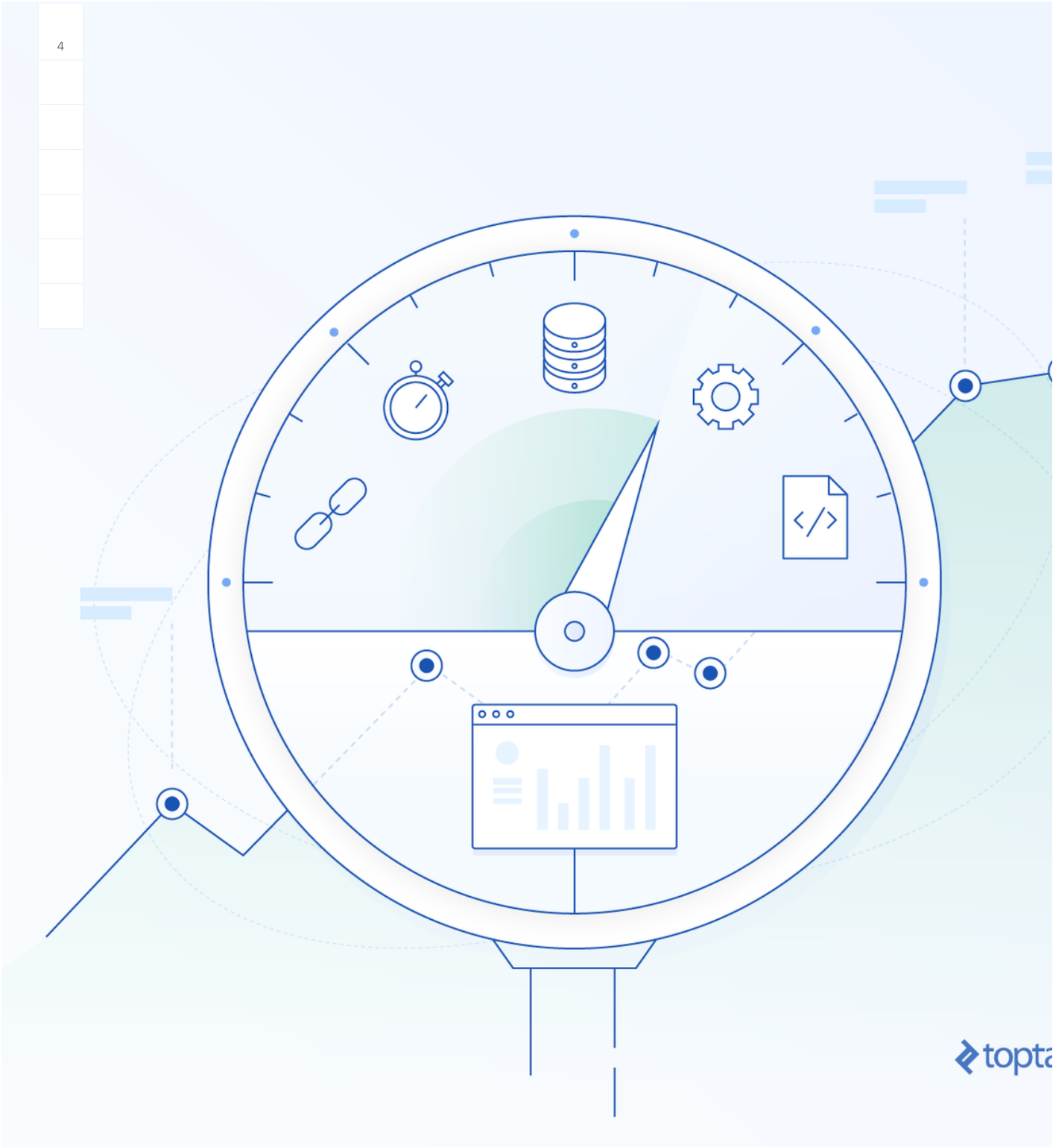
原

作

翻

摘要：通过一个简单的实例一步一步引导读者对其进行全方位的性能优化。以下是译文。

唐纳纳德·克特（Donald Knuth）曾经说过：“不成熟的优化方案是万恶之源。”然而，任何一个承受高负载的成熟项目都不可避免地需要进行优化。在本文中，我们将谈谈优化Web项目代码的五种常用方法。虽然本文是以Django为例，但其他框架和语言的优化原则也是类似的。通过使用这些优化方法，文中例程的执行时间从原来的77秒减少到了3.7秒。



本文用到的例程是从一个我曾经使用过的真实项目改编而来的，是性能优化技巧的典范。如果你想自己尝试着进行优化，可以在[GitHub](#)上获取优化前并跟着下文做相应的修改。我使用的是Python 2，因为一些第三方软件包还不支持Python 3。

示例代码介绍

这个Web项目只是简单地跟踪每个地区的房产价格。因此，只有两种模型：

```
1 # houses/models.py
2 from utils.hash import Hasher
3
4
5 class HashableModel(models.Model):
6     """Provide a hash property for models."""
7     class Meta:
```

```

8         abstract = True
9
10    @property
11    def hash(self):
12        4         return Hasher.from_model(self)
13
14
15    ss Country(HashableModel):
16        """Represent a country in which the house is positioned."""
17        name = models.CharField(max_length=30)
18
19        def __unicode__(self):
20            return self.name
21
22
23    ss House(HashableModel):
24        """Represent a house with its characteristics."""
25        # Relations
26        country = models.ForeignKey(Country, related_name='houses')
27
28        # Attributes
29        address = models.CharField(max_length=255)
30        sq_meters = models.PositiveIntegerField()
31        kitchen_sq_meters = models.PositiveSmallIntegerField()
32        nr_bedrooms = models.PositiveSmallIntegerField()
33        nr_bathrooms = models.PositiveSmallIntegerField()
34        nr_floors = models.PositiveSmallIntegerField(default=1)
35        year_built = models.PositiveIntegerField(null=True, blank=True)
36        house_color_outside = models.CharField(max_length=20)
37        distance_to_nearest_kindergarten = models.PositiveIntegerField(null=True, blank=True)
38        distance_to_nearest_school = models.PositiveIntegerField(null=True, blank=True)
39        distance_to_nearest_hospital = models.PositiveIntegerField(null=True, blank=True)
40        has_cellar = models.BooleanField(default=False)
41        has_pool = models.BooleanField(default=False)
42        has_garage = models.BooleanField(default=False)
43        price = models.PositiveIntegerField()
44
45        def __unicode__(self):
46            return '{} {}'.format(self.country, self.address)

```

抽象类 `HashableModel` 提供了一个继承自模型并包含 `hash` 属性的模型，这个属性包含了实例的主键和模型的内容类型。这能够隐藏像实例ID这样的敏感列进行代替。如果项目中有多个模型，而且需要在一个集中的地方对模型进行解码并要对不同类的不同模型实例进行处理时，这可能会非常有用。请文的这个小项目，即使不用散列也照样可以处理，但使用散列有助于展示一些优化技巧。

这是 `Hasher` 类：

```

1  # utils/hash.py
2  import basehash
3
4
5  class Hasher(object):
6      @classmethod
7      def from_model(cls, obj, klass=None):
8          if obj.pk is None:
9              return None
10         return cls.make_hash(obj.pk, klass if klass is not None else obj)
11
12     @classmethod
13     def make_hash(cls, object_pk, klass):
14         base36 = basehash.base36()
15         content_type = ContentType.objects.get_for_model(klass, for_concrete_model=False)
16         return base36.hash('%(contenttype_pk)03d%(object_pk)06d' % {
17             'contenttype_pk': content_type.pk,
18             'object_pk': object_pk
19         })
20
21     @classmethod
22     def parse_hash(cls, obj_hash):
23         base36 = basehash.base36()
24         unhashed = '%09d' % base36.unhash(obj_hash)

```

```

25         contenttype_pk = int(unhashed[:-6])
26         object_pk = int(unhashed[-6:])
27         return contenttype_pk, object_pk
28
29     @classmethod
30     def to_object_pk(cls, obj_hash):
31         return cls.parse_hash(obj_hash)[1]

```

由于我们希望通过API来提供这些数据，所以我们安装了Django REST框架并定义以下序列化器和视图：

```

1 # houses/serializers.py
2 from rest_framework import serializers
3
4
5 class HouseSerializer(serializers.ModelSerializer):
6     """Serialize a `houses.House` instance."""
7
8     id = serializers.ReadOnlyField(source="hash")
9     country = serializers.ReadOnlyField(source="country.hash")
10
11     class Meta:
12         model = House
13         fields = (
14             'id',
15             'address',
16             'country',
17             'sq_meters',
18             'price'
19         )
20
21
22 # houses/views.py
23 from rest_framework import generics
24
25 class HouseListAPIView(generics.ListAPIView):
26     model = House
27     serializer_class = HouseSerializer
28     country = None
29
30     def get_queryset(self):
31         country = get_object_or_404(Country, pk=self.country)
32         queryset = self.model.objects.filter(country=country)
33         return queryset
34
35     def list(self, request, *args, **kwargs):
36         # Skipping validation code for brevity
37         country = self.request.GET.get("country")
38         self.country = Hasher.to_object_pk(country)
39         queryset = self.get_queryset()
40
41         serializer = self.serializer_class(queryset, many=True)
42
43         return Response(serializer.data)

```

现在，我们将用一些数据来填充数据库（使用 **factory-boy** 生成10万个房屋的实例：一个地区5万个，另一个4万个，第三个1万个），并准备测试应用程序。

性能优化其实就是测量

在一个项目中我们需要测量下面这几个方面：

- 执行时间
- 代码的行数
- 函数调用次数
- 分配的内存
- 其他

但是，并不是所有这些都要用来度量项目的执行情况。一般来说，有两个指标比较重要：执行多长时间、需要多少内存。

在Web项目中，**响应时间**（服务器接收由某个用户的操作产生的请求，处理该请求并返回结果所需的总的时间）通常是最重要的指标，因为过长的响应时间会导致用户厌倦等待，并切换到浏览器中的另一个选项卡页面。

在编程中，分析项目的性能被称为**profiling**。为了分析API的性能，我们将使用**Silk**包。在安装完这个包，并调用 `/api/v1/houses/?country=5T22RI` 后，我们将看到以下结果：

```
1 200 GET
2 /api/v1/houses/
3
4 92ms overall
5 4 54ms on queries
6 34 queries
```

整体时间为77秒，其中16秒用于查询数据库，总共有5万次查询。这几个数字很大，提升空间也有很大，所以，我们开始吧。

1.1 数据库查询

性能提升常见的技巧之一是对数据库查询进行优化，本案例也不例外。同时，还可以对查询做多次优化来减小响应时间。

1.1 提供所有数据

仔细分析这5万次查询查的是什么：都是对 `houses_country` 表的查询：

```
1 200 GET
2 /api/v1/houses/
3
4 77292ms overall
5 15854ms on queries
6 50004 queries
```

时间戳	表名	联合	执行时间（毫秒）
+0:01:15.874374	"houses_country"	0	0.176
+0:01:15.873304	"houses_country"	0	0.218
+0:01:15.872225	"houses_country"	0	0.218
+0:01:15.871155	"houses_country"	0	0.198
+0:01:15.870099	"houses_country"	0	0.173
+0:01:15.869050	"houses_country"	0	0.197
+0:01:15.867877	"houses_country"	0	0.221
+0:01:15.866807	"houses_country"	0	0.203
+0:01:15.865646	"houses_country"	0	0.211
+0:01:15.864562	"houses_country"	0	0.209
+0:01:15.863511	"houses_country"	0	0.181
+0:01:15.862435	"houses_country"	0	0.228
+0:01:15.861413	"houses_country"	0	0.174

这个问题的根源是，Django中的查询是惰性的。这意味着在你真正需要获取数据之前它不会访问数据库。同时，它只获取你指定的数据，如果需要其他数据，则需要另外发出请求。

这正是本例程所遇到的情况。当通过 `House.objects.filter(country=country)` 来获得查询集时，Django将获取特定地区的所有房屋。但是，在序列化一时，`HouseSerializer` 需要房子的 `country` 实例来计算序列化器的 `country` 字段。由于地区数据不在查询集中，所以django需要提出额外的请求来获取这。查询集中的每一个房子都是如此，因此，总共是五万次。

当然，解决方案非常简单。为了提取所有需要的序列化数据，你可以在查询集上使用 `select_related()`。因此，`get_queryset` 函数将如下所示：

```
1 def get_queryset(self):
2     country = get_object_or_404(Country, pk=self.country)
3     queryset = self.model.objects.filter(country=country).select_related('country')
4     return queryset
```

我们来看看这对性能有何影响：

```
1 200 GET
2 /api/v1/houses/
3
```

```

4 35979ms overall
5 102ms on queries
6 4 queries

```

总体时间降至36秒，在数据库中花费的时间约为100ms，只有4个查询！这是个好消息，但我们可以做得更多。

1.2 提供相关的数据

默认情况下，Django会从数据库中提取所有字段。但是，当表有很多列很多行的时候，告诉Django提取哪些特定的字段就非常有意义了，这样就不会根本性的信息。在本案例中，我们只需要5个字段来进行序列化，虽然表中有17个字段。明确指定从数据库中提取哪些字段是很有意义的，可以节省时间。

Django使用`defer()`和`only()`这两个查询方法来实现这一点。第一个用于指定哪些字段不要加载，第二个用于指定只加载哪些字段。

```

1 get_queryset(self):
2     country = get_object_or_404(Country, pk=self.country)
3     queryset = self.model.objects.filter(country=country)\
4         .select_related('country')\
5         .only('id', 'address', 'country', 'sq_meters', 'price')
6     return queryset

```

这减少了一半的查询时间，非常不错。总体时间也略有下降，但还有更多提升空间。

```

1 200 GET
2 /api/v1/houses/
3
4 3311ms overall
5 52ms on queries
6 4 queries

```

2. 代码优化

你不能无限地优化数据库查询，并且上面的结果也证明了这一点。即使把查询时间减少到0，我们仍然会面对需要等待半分钟才能得到应答这个现实转移到另一个优化级别上来了，那就是：*业务逻辑*。

2.1 简化代码

有时，第三方软件包对于简单的任务来说有着太大的开销。本文例程中返回的序列化的房子实例正说明了这一点。

Django REST框架非常棒，包含了很多有用的功能。但是，现在的主要目标是缩短响应时间，所以该框架是优化的候选对象，尤其是我们要使用的序列化功能非常的简单。

为此，我们来编写一个自定义的序列化器。为了方便起见，我们将用一个静态方法来完成这项工作。

```

1 # houses/serializers.py
2 class HousePlainSerializer(object):
3     """
4     Serializes a House queryset consisting of dicts with
5     the following keys: 'id', 'address', 'country',
6     'sq_meters', 'price'.
7     """
8
9     @staticmethod
10    def serialize_data(queryset):
11        """
12        Return a list of hashed objects from the given queryset.
13        """
14        return [
15            {
16                'id': Hasher.from_pk_and_class(entry['id'], House),
17                'address': entry['address'],
18                'country': Hasher.from_pk_and_class(entry['country'], Country),
19                'sq_meters': entry['sq_meters'],
20                'price': entry['price']
21            } for entry in queryset
22        ]
23
24
25 # houses/views.py

```

```

26 class HouseListAPIView(ListAPIView):
27     model = House
28     serializer_class = HouseSerializer
29     plain_serializer_class = HousePlainSerializer # <-- added custom serializer
30     4 country = None
31
32     def get_queryset(self):
33         country = get_object_or_404(Country, pk=self.country)
34         queryset = self.model.objects.filter(country=country)
35         return queryset
36
37     def list(self, request, *args, **kwargs):
38         # Skipping validation code for brevity
39         country = self.request.GET.get("country")
40         self.country = Hasher.to_object_pk(country)
41         queryset = self.get_queryset()
42
43         data = self.plain_serializer_class.serialize_data(queryset) # <-- serialize
44
45         return Response(data)

```



```

1 200 GET
2 /api/v1/houses/
3
4 17312ms overall
5 38ms on queries
6 4 queries

```

现在看起来好多了，由于没有使用DRF序列化代码，所以响应时间几乎减少了一半。

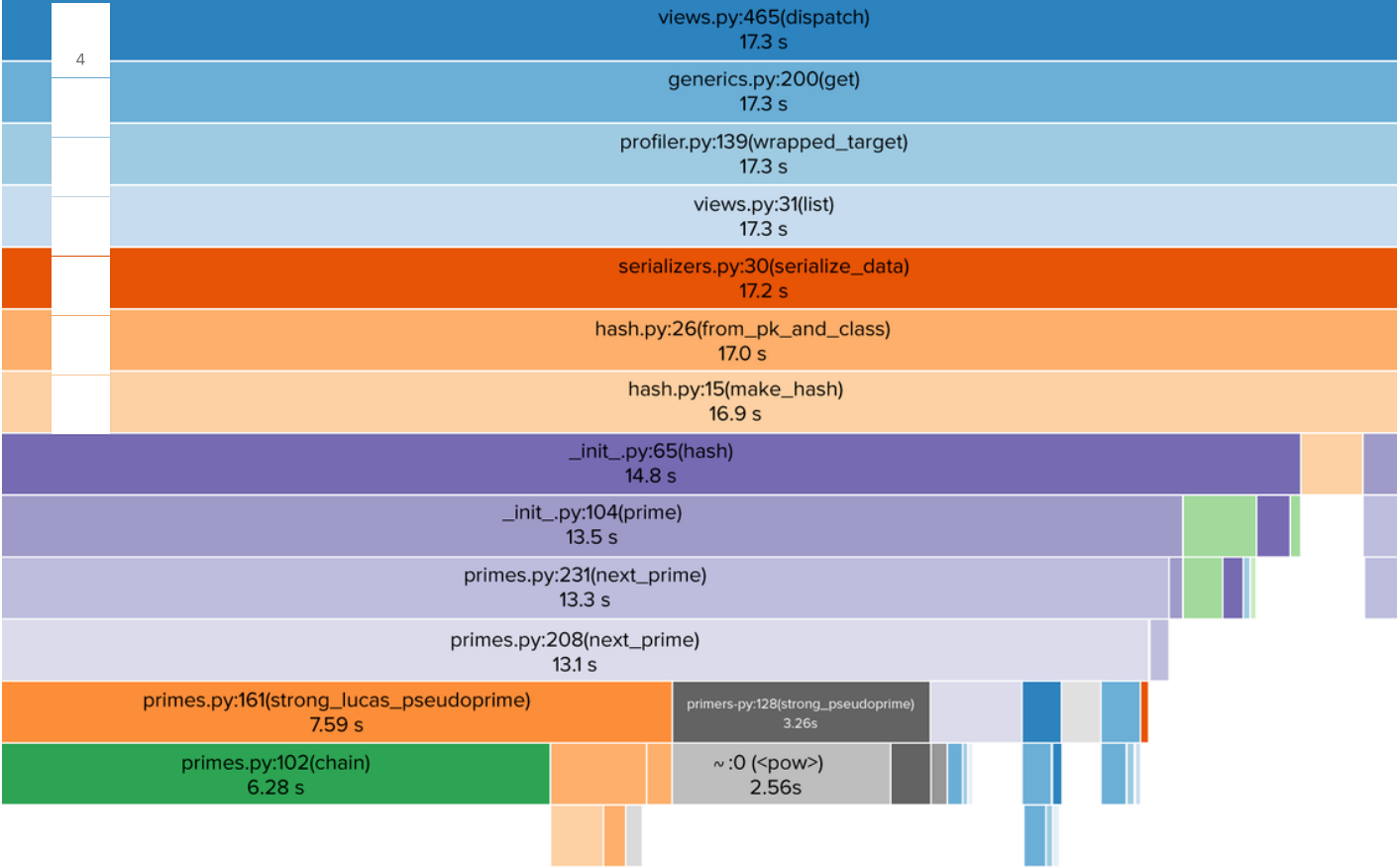
另外还有一个结果：在请求/响应周期内完成的总的函数调用次数从15,859,427次（上面1.2节的请求次数）减少到了9,257,469次。这意味着大约有三调用都是由Django REST Framework产生的。

2.2 更新或替代第三方软件包

上述几个优化技巧是最常见的，无需深入地分析和思考就可以做到。然而，17秒的响应时间仍然感觉很长。要减少这个时间，需要更深入地了解代码生了什么。换句话说，需要分析一下代码。

你可以自己使用Python内置的分析器来进行分析，也可以使用一些第三方软件包。由于我们已经使用了 `silk`，它可以分析代码并生成一个二进制的：此，我们可以做进一步的可视化分析。有好几个可视化软件包可以将二进制文件转换为一些友好的可视化视图。本文将使用 `snakeviz`。

这是上文一个请求的二进制分析文件的可视化图表：



从上到下是调用堆栈，显示了文件名、函数名及其行号，以及该方法花费的时间。可以很容易地看出，时间大部分都用在计算散列上（紫罗兰色的 `_init_.py` 矩形）。

目前，这是代码的主要性能瓶颈，但同时，这不是我们自己写的代码，而是用的第三方包。

在这种情况下，我们可以做的事情将非常有限：

- 检查包的最新版本（希望能有更好的性能）。
- 寻找另一个能够满足我们需求的软件包。
- 我们自己写代码，并且性能优于目前使用的软件包。

幸运的是，我们找到了一个更新版本的 `basehash` 包。原代码使用的是v.2.1.0，而新的是v.3.0.4。

当查看v.3的发行说明时，这一句话看起来令人充满希望：

“使用素数算法进行大规模的优化。”

让我们来看一下！

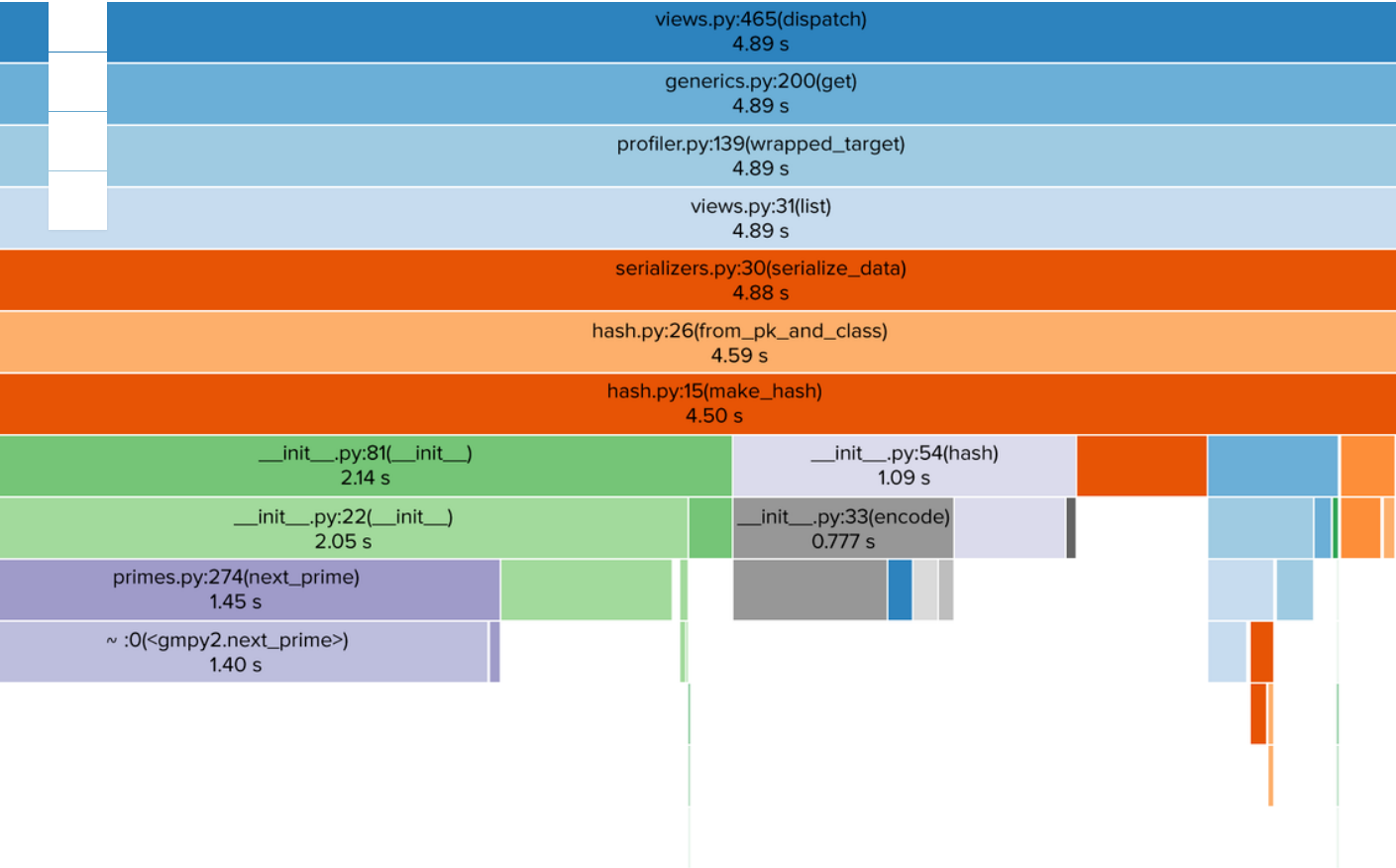
```
1 pip install -U basehash gmpy2

1 200 GET
2 /api/v1/houses/
3
4 7738ms overall
5 59ms on queries
6 4 queries
```


响应时间从17秒缩短到了8秒以内。太棒了！但还有一件事我们应该来看看。

2.3 代码

到目前为止，我们已经改进了查询、用自己特定的函数取代了第三方复杂而又泛型的代码、更新了第三方包，但是我们还是保留了原有的代码。但有时进行代码的重构可能会带来意想不到的结果。但是，为此我们需要再次分析运行结果。



仔细看一下，你可以看到散列仍然是一个问题（毫不奇怪，这是我们对数据做的唯一的事情），虽然我们确实朝这个方向改进了，但这个绿色的矩形 `__init__.py` 花了2.14秒的时间，同时伴随着灰色的 `__init__.py:54(hash)`。这意味着初始化工作需要很长的时间。

我们来看看 `basehash` 包的源代码。

```
1 # basehash/__init__.py
2
3 # Initialization of `base36` class initializes the parent, `base` class.
4 class base36(base):
5     def __init__(self, length=HASH_LENGTH, generator=GENERATOR):
6         super(base36, self).__init__(BASE36, length, generator)
7
8
9 class base(object):
10     def __init__(self, alphabet, length=HASH_LENGTH, generator=GENERATOR):
11         if len(set(alphabet)) != len(alphabet):
12             raise ValueError('Supplied alphabet cannot contain duplicates.')
13
14         self.alphabet = tuple(alphabet)
15         self.base = len(alphabet)
16         self.length = length
17         self.generator = generator
18         self.maximum = self.base ** self.length - 1
19         self.prime = next_prime(int((self.maximum + 1) * self.generator)) # `next_prime` call on each initialized instance
```

正如你所看到的，一个 `base` 实例的初始化需要调用 `next_prime` 函数，这是太重了，我们可以在上面的可视化图表中看到左下角的矩形。

我们再看 `Hash` 类：

```

1  class Hasher(object):
2      @classmethod
3      def from_model(cls, obj, klass=None):
4          if obj.pk is None:
5              return None
6          return cls.make_hash(obj.pk, klass if klass is not None else obj)
7
8      @classmethod
9      def make_hash(cls, object_pk, klass):
10         base36 = basehash.base36() # <-- initializing on each method call
11         content_type = ContentType.objects.get_for_model(klass, for_concrete_model=False)
12         return base36.hash('%(contenttype_pk)03d%(object_pk)06d' % {
13             'contenttype_pk': content_type.pk,
14             'object_pk': object_pk
15         })
16
17     @classmethod
18     def parse_hash(cls, obj_hash):
19         base36 = basehash.base36() # <-- initializing on each method call
20         unhashed = '%09d' % base36.unhash(obj_hash)
21         contenttype_pk = int(unhashed[:-6])
22         object_pk = int(unhashed[-6:])
23         return contenttype_pk, object_pk
24
25     @classmethod
26     def to_object_pk(cls, obj_hash):
27         return cls.parse_hash(obj_hash)[1]
```

正如你所看到的，我已经标记了这两个方法初始化 `base36` 实例的方法，这并不是真正需要的。

由于散列是一个确定性的过程，这意味着对于一个给定的输入值，它必须始终生成相同的散列值，因此，我们可以把它作为类的一个属性。让我们来看行：

```

1  class Hasher(object):
2      base36 = basehash.base36() # <-- initialize hasher only once
3
4      @classmethod
5      def from_model(cls, obj, klass=None):
6          if obj.pk is None:
7              return None
8          return cls.make_hash(obj.pk, klass if klass is not None else obj)
9
10     @classmethod
11     def make_hash(cls, object_pk, klass):
12         content_type = ContentType.objects.get_for_model(klass, for_concrete_model=False)
13         return cls.base36.hash('%(contenttype_pk)03d%(object_pk)06d' % {
14             'contenttype_pk': content_type.pk,
15             'object_pk': object_pk
16         })
17
18     @classmethod
19     def parse_hash(cls, obj_hash):
20         unhashed = '%09d' % cls.base36.unhash(obj_hash)
21         contenttype_pk = int(unhashed[:-6])
22         object_pk = int(unhashed[-6:])
23         return contenttype_pk, object_pk
24
25     @classmethod
26     def to_object_pk(cls, obj_hash):
27         return cls.parse_hash(obj_hash)[1]
```

```

1  **200 GET**
2
3  /api/v1/houses/
4
5  3766ms overall
```

```
6 38ms on queries
7 4 queries
```

最后4是在4秒钟之内，比我们一开始的时间要小得多。对响应时间的进一步优化可以通过使用缓存来实现，但是我不会在这篇文章中介绍这个。

结论

性能优化是一个分析和发现的过程。没有哪个硬性规定能适用于所有情况，因为每个项目都有自己的流程和瓶颈。然而，你应该做的第一件事是分析你这样做的例子中，我可以将响应时间从77秒缩短到3.7秒，那么对于一个庞大的项目来说，就会有更大的优化潜力。

PS：给大家推荐一个数据库在线峰会，嘉宾来自阿里巴巴、腾讯、微博、网易等多家企业的数据库专家及高校研究学者，将围绕Kubernetes、MongoDB、Redis等热点数据库技术展开。峰会详情页：http://edu.csdn.net/huiyiCourse/series_detail/74



文章标签：python django 性能测试 ▼查看关于本篇文章更多信息

想对作者说点什么

提高django性能 444
性能优化是一件困难的事情，但是也不常常如此：下面4步将能够轻松的提高你的网站的性能，它们非常简单你...

抒发一下这些天用django做web项目的一些体会 4085
最近接触了一段时间的python，觉得python写脚本还是挺方便的，做一个简单的桌面应用也很nice，但是随着深...

Django 性能测试——一个现实世界的例子 - CSDN博客
“墙”页面的性能令人失望,因此我们开始进行优化。我们所做第一件事情是分析代码...测了一下django、flask、bo...

优化Django ORM中的性能问题 - CSDN博客
如果 prefetch太复杂了,这时候就要在代码的整洁清晰和应用性能之间做一个取舍了...了解ORM 是怎么缓存数据...

Django框架效率问题的解决方法和总... 1286
由于项目的需要，学习了Django框架，Django框架的MTV很清晰，通过MTV能够很好地了解Django框架的内部...

优化Django ORM中的性能问题 3133
原文地址 Solving Performance Problems in the Django ORM Django是个好工具，使用的很广泛。在应用比较...

如何进行Django单元测试? - CSDN博客
Django的单元测试使用python的unittest模块,这个模块使用基于类的方法来定义测试。类名为django.test.TestCa...

cpu性能测试代码 - CSDN博客
CPU-Pi -- 基于 Java 的 CPU性能测试工具 及源代码 qq_25670227 05-18 73...Django 5篇 dsst 8篇 Eigen 8篇 io...