

django in the real world

yes! it scales!... YAY!

Israel Fermin Montilla

Software Engineer @ dubizzle

December 14, 2017

from iferminm import more_data

- ▶ Software Engineer @ dubizzle
- ▶ Venezuelan living in Dubai, UAE
- ▶ T: @iferminm
- ▶ blog: <http://iffm.me>

What will we see in this talk?

- ▶ Pareto Principle
- ▶ The simple django project
- ▶ Measuring
- ▶ Common bottlenecks

Basic concepts: Pareto principle

Basic concepts: Pareto principle

The Pareto principle states that, for many events, roughly 80% of the effects come from 20% of the causes

–Wikipedia

Basic concepts: Pareto principle

The Pareto principle states that, for many events, roughly 80% of the effects come from 20% of the causes

–Wikipedia

For example: 20% of the code produces 80% of the bugs.

Basic concepts: Pareto principle

- ▶ Premature optimization is bad

Basic concepts: Pareto principle

- ▶ Premature optimization is bad
- ▶ Optimization without measuring is bad

Basic concepts: Pareto principle

- ▶ Premature optimization is bad
- ▶ Optimization without measuring is bad
- ▶ Unprioritized optimization is bad

Initial django project in production

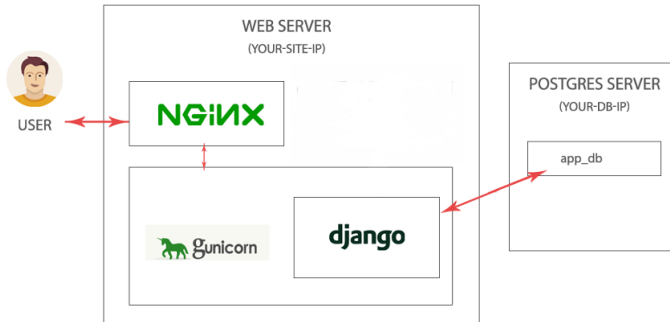


Figure: Basic django project production setup

Profile first



django-debug-toolbar

The screenshot displays the Django Debug Toolbar interface. The main panel is titled "SQL queries from 1 connection" and shows a list of five SQL queries with their execution times and actions. The queries are:

Query	Timeline	Time (ms)	Action
<code>SELECT ... FROM "django_session" WHERE ("django_session"."session_key" = 'n30ikmlcaz0e7j1hte50a4x01z5hwt6c' AND "django_session"."expire_date" > '2013-12-21 03:26:09.232862')</code>	[Timeline bar]	1.71	Sel Expl
<code>SELECT ... FROM "auth_user" WHERE "auth_user"."id" = 1</code>	[Timeline bar]	1.01	Sel Expl
<code>SELECT ... FROM "auth_group"</code>	[Timeline bar]	0.43	Sel Expl
<code>SELECT ... FROM "auth_user"</code>	[Timeline bar]	0.68	Sel Expl
<code>SELECT ... FROM "auth_user" ORDER BY "auth_user"."username" ASC, "auth_user"."id" DESC</code>	[Timeline bar]	1.19	Sel Expl

The sidebar on the right contains several sections:

- Hide »** (with sub-options: password, log out)
- Versions** (checked): DJANGO 1.6
- Time** (checked): CPU: 119.44ms (126.50ms)
- Settings** (checked): of superuser status
- Headers** (checked)
- Request** (checked)
- SQL** (checked): 5 QUERIES IN 5.02MS
- Templates** (checked)
- Static files** (checked): 10 FILES USED

Figure: debug_toolbar in action

cProfile + snakeviz

Search:

ncalls ↕	tottime ▼	percall ↕	cumtime ↕	percall ↕	filename:lineno(function) ↕
1	0.000421	0.000421	0.000421	0.000421	~:0(<built-in method listdir>)
1	0.000104	0.000104	0.000202	0.000202	functools.py:441(wrapper)
1	7.9e-05	7.9e-05	0.000294	0.000294	fnmatch.py:48(filter)
1	6.7e-05	6.7e-05	8e-05	8e-05	functools.py:342(_make_key)
1	4.4e-05	4.4e-05	0.00079	0.00079	glob.py:61(glob1)

Figure: snakeviz list view

cProfile + snakeviz

Name:

filter

Cumulative Time:

0.000294 s (31.78 %)

File:

fnmatch.py

Line:

48

Directory:

/Users/jiffyclub/miniconda3/en
vs/snakevizdev/lib/python3.4/

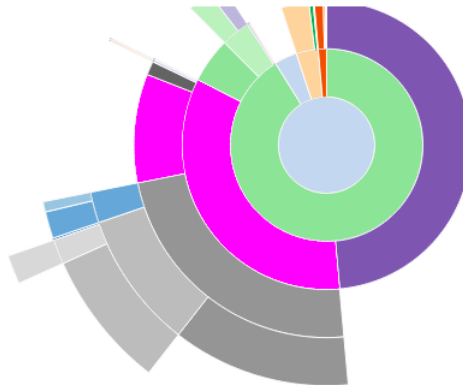


Figure: snakeviz sunburst diagram

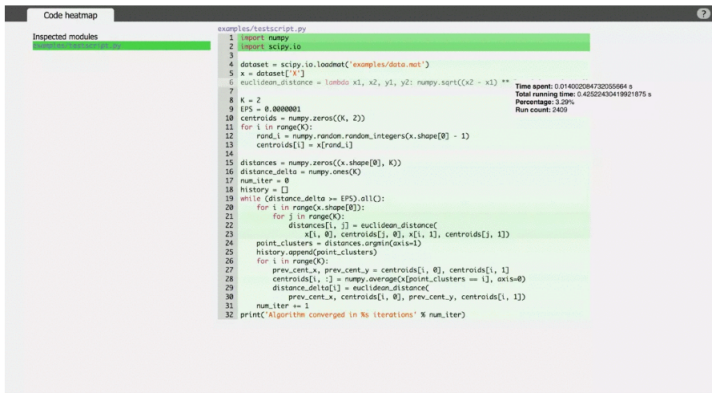


Figure: vprof code heatmap

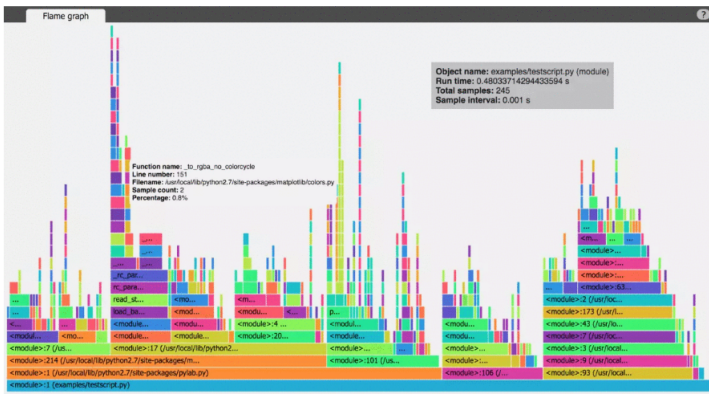


Figure: vprof flame diagram

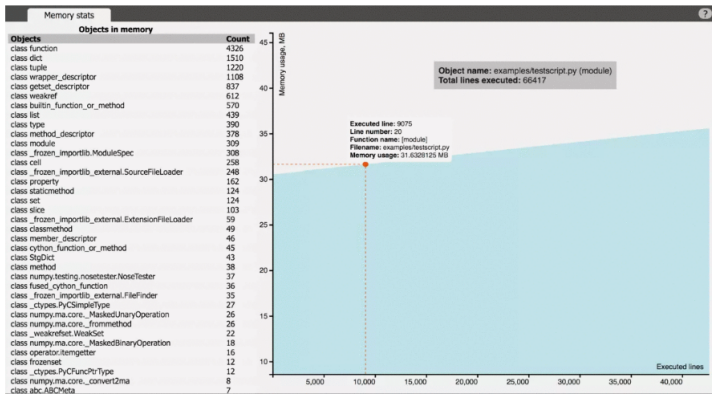


Figure: vprof memory profiler

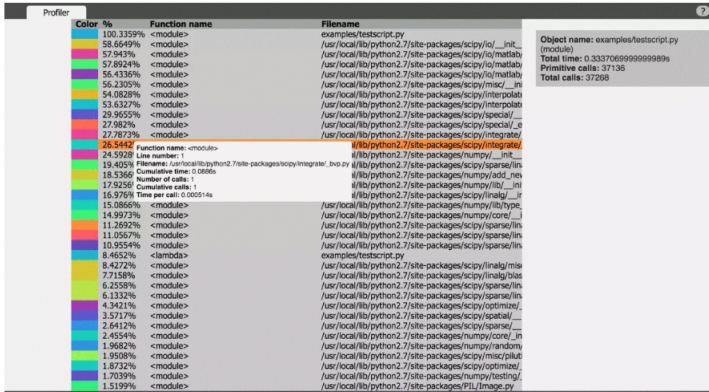
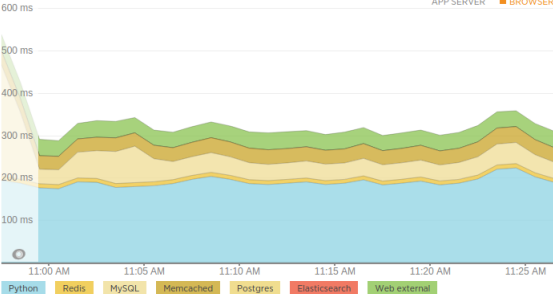
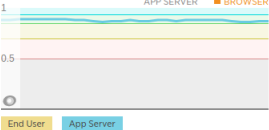


Figure: vprof profiler

Web transactions time



Apdex score



Throughput

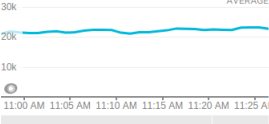


Figure: Part of newrelic's main dashboard

Transactions

App server time

`/classified.views:mylistings` 4.02 sec

Transaction traces: 36.9 s 23.2 s 17.7 s

`/kombi.views:ListingView.inner` 770 ms

Transaction traces: 4.7 s 4.1 s 3.8 s

`/generic_api.slayer.views:UserAPIView` 578 ms

Transaction traces: 3.1 s 2.9 s 2.8 s

`/classified.views:category_list` 557 ms

Transaction traces: 3 s 2.6 s 2.5 s

`/classified.views:classified_detail` 509 ms

Transaction traces: 8 s 3.7 s 3.3 s

Error rate

0.0139 %



Breakout each metric by host

Figure: part of newrelic's main dashboard

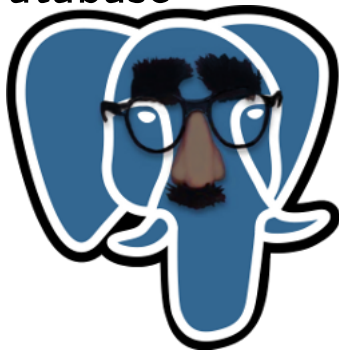
Breakdown table

Category	Segment	% Time	Avg calls (per txn)	Avg time (ms)
Function	kombi.views:ListingView.inner	63.1	1.0	491
Database	Memcached get	4.3	36.7	33.4
Database	Memcached set	1.1	4.76	8.87
Function	MySQLdb:Connect	0.9	1.98	6.86
Database	MySQL classified_classified_au select	0.8	2.09	6.06
Database	MySQL classified_propertyforrent_rs select	0.7	1.91	5.62
Database	MySQL classified_classified_fu select	0.3	1.95	2.26
Database	MySQL classified_propertyforsale_rs select	0.3	1.86	2.54

[Show all segments →](#)

Figure: Inside a web transaction in newrelic

Database



Reduce query counts

```
1 subs = Subscription.objects.filter(user_id=user.pk)
2 for s in subs:
3     packages.append(s.package.name)
```

Reduce query counts

```
1 subs = Subscription.objects.filter(user_id=user.pk)
2 for s in subs:
3     packages.append(s.package.name)
```

N hits to the database

Reduce query counts

```
1 subs = Subscription.objects.filter(user_id=user.pk)
2 for s in subs:
3     packages.append(s.package.name)
```

N hits to the database

```
1 subs = Subscription.objects.filter(
2     user_id=user.pk
3 ).select_related('package')
```

Reduce query counts

```
1 subs = Subscription.objects.filter(user_id=user.pk)
2 for s in subs:
3     packages.append(s.package.name)
```

N hits to the database

```
1 subs = Subscription.objects.filter(
2     user_id=user.pk
3 ).select_related('package')
```

Will join the table and return it in one hit

Reduce query counts

- ▶ `select_related`
- ▶ `prefetch_related`

Reduce query counts

- ▶ Use it wisely and measure

```
1 user = User.objects.select_related(  
2     'sodas'  
3 ).get(pk=request.data['user_id'])  
4  
5 # No additional query  
6 user.sodas.all()
```

Reduce query counts

► Use it wisely and measure

```
1 user = User.objects.select_related(  
2     'sodas'  
3 ).get(pk=request.data['user_id'])  
4  
5 # No additional query  
6 user.sodas.all()
```

```
1 # Triggers an additional query  
2 user.sodas.filter(name='pepsi')  
3  
4 # Sometimes it's better to use the cached result  
5 # and filter in memory  
6 [s for s in user.sodas.all() if s.name == 'pepsi']
```

Reduce query counts

Use the Prefetch object!

```
1 # A product has many subscriptions and
2 # a subscription can have many products
3
4 queryset = Subscription.objects.filter(status=expired).select_related('credits')
5 prefetch = Prefetch('subscriptions', queryset=queryset)
6
7 products = Product.objects.prefetch_related(prefetch).filter(section='jobs')
```

Reduce query time

Reduce query time

- ▶ Indexing

Reduce query time

► Indexing

```
1 class UserProfile(models.Model):  
2     user = models.ForeignKey('auth.user')  
3     dob = models.DateField(db_index=True)  
4     external_id = models.IntegerField(db_index=True)
```

Reduce query time

► Indexing

```
1 class UserProfile(models.Model):  
2     user = models.ForeignKey('auth.user')  
3     dob = models.DateField(db_index=True)  
4     external_id = models.IntegerField(db_index=True)
```

Note: Your DBMS updates your indices in *write time* (*INSERT and UPDATE*)

Some notes on indexing

- ▶ You need to measure before you do it. Run EXPLAIN on the query (Seq scan)
- ▶ Index by workload
- ▶ If you filter on multiple columns use index_together Meta option
- ▶ Check if the index is used before you push it. Run EXPLAIN again

Expensive JOINS

Sometimes you might want to separate them into two different queries.

```
1 # You may want to see the credit spending behavior of your users
2 Credit.objects.filter(
3     subscription__pkg__type='motors'
4 ).select_related('resource')
5
6 # Sometimes two queries might perform better
7 subs_ids = Subscription.objects.filter(
8     pkg__type='motors'
9 ).values_list('id', flat=True)
10
11 Credit.objects.filter(
12     subscription_id__in=subs_ids
13 ).select_related('resource')
```

Expensive JOINS

Sometimes you might want to separate them into two different queries.

```
1 # You may want to see the credit spending behavior of your users
2 Credit.objects.filter(
3     subscription__pkg__type='motors'
4 ).select_related('resource')
5
6 # Sometimes two queries might perform better
7 subs_ids = Subscription.objects.filter(
8     pkg__type='motors'
9 ).values_list('id', flat=True)
10
11 Credit.objects.filter(
12     subscription_id__in=subs_ids
13 ).select_related('resource')
```

ALWAYS MEASURE

Avoid whole table COUNT() queries

After some point, having exact numbers is not important

```
1 PropertyForRent.objects.count()
```

Avoid whole table COUNT() queries

After some point, having exact numbers is not important

```
1 PropertyForRent.objects.count()
```

You can instead do a raw SQL query

```
1 # Postgres
2 SELECT reltuples FROM pg_class
3     WHERE relname = 'property_for_rent'
4
5 # MySQL
6 SELECT table_rows FROM information_schema.tables
7     WHERE table_schema = DATABASE()
8     AND table_name = 'property_for_rent'
```

Avoid whole table COUNT() queries

After some point, having exact numbers is not important

```
1 PropertyForRent.objects.count()
```

You can instead do a raw SQL query

```
1 # Postgres
2 SELECT reltuples FROM pg_class
3     WHERE relname = 'property_for_rent'
4
5 # MySQL
6 SELECT table_rows FROM information_schema.tables
7     WHERE table_schema = DATABASE()
8     AND table_name = 'property_for_rent'
```

This could reduce up to 90% response time

Use persistent connections

```
1 DATABASES = {  
2     'default': {  
3         'ENGINE': 'django.db.backends.postgresql_psycopg2',  
4         'NAME': os.getenv('DATABASE_NAME', 'slayer'),  
5         'USER': os.getenv('DATABASE_USER', None),  
6         'PASSWORD': os.getenv('DATABASE_PASSWORD', None),  
7         'PORT': os.getenv('DATABASE_PORT', '3306'),  
8         'HOST': os.getenv('DATABASE_HOST', 'localhost'),  
9         'CONN_MAX_AGE': int(os.getenv('DATABASE_CONNECTION_MAX_AGE', '0'))  
10     }  
11 }
```

Know your ORM

- ▶ Read the full ORM docs at least once
- ▶ Use F expressions to reference values within the queryset
- ▶ Use Q expressions for advanced filters
- ▶ Explore the aggregation framework
- ▶ Use `values()`, `values_list()`, `only()` and `defer()` when the results are too big

Denormalize

- ▶ Evaluate huge joins
- ▶ Don't use Generic Relations

Denormalize

- ▶ Evaluate huge joins
- ▶ Don't use Generic Relations

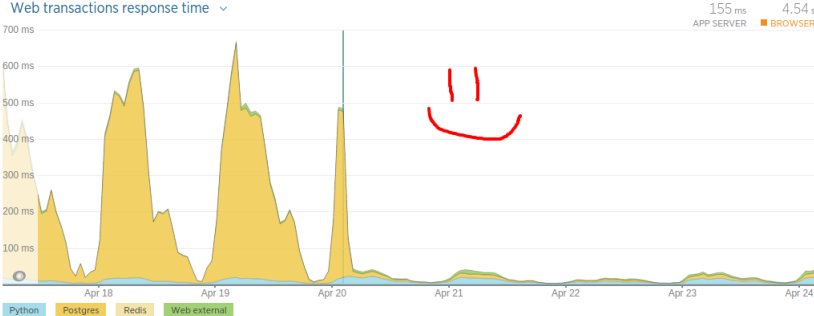


Figure: Response time reduction after denormalizing a Generic Relation

Query caching

- ▶ johny-cache
- ▶ django-cache-machine

Templates

Russian Doll Caching



Russian Doll Caching

```
1 {% cache MIDDLE_TTL "ads" request.GET.page %}
2     {% include "sections/property/postheader.html" %}
3     <div class="ads-list">
4         {% for ad in ads %}
5             {% cache LONG_TTL "ad_description" ad.id ad.last_updated %}
6                 {% include "sections/property/ad_teaser.html" %}
7             {% endcache %}
8         {% endfor %}
9     {% endcache %}
```


Further Optimization

Further optimization

- ▶ Minimize your CSS and JS (django-compressor, webassets or django-pipeline)
- ▶ Optimize your static images
- ▶ Optimize user uploaded images
- ▶ Serve your media and static content from a CDN
- ▶ Do slow work later... (celery or python-rq)
- ▶ Use slave replicas for *read operations (and database routers)*

DESPACITO

Thank you!



**** We're hiring ****
**** ifermin@gmail.com ****