



nyt.tgz

1998/06/15/1024820.xml  
1998/06/15/1024821.xml

```
<?xml version="1.0" encoding="UTF-8" ?>
.....
.....
.....
<p>J P Hayes wins Buick Classic golf
tournament; photos (M) </p>
.....
.....
```

→  
Cleaning  
and  
xml processing



nyt\_cleaned.tgz

1024/1024820  
1024/1024821

```
.....
.....
.....
J P Hayes wins Buick Classic golf
tournament; photos (M)
.....
.....
```

→  
Corpus  
Preprocessor



nyt\_features.bin