# IMDb

At Fast.ai we have introduced a new module called fastai.text which replaces the torchtext library that was used in our 2018 dl1 course. The fastai.text module also supersedes the fastai.nlp library but retains many of the key functions.

```
In [1]: from fastai.text import *
        import html
```

The Fastai.text module introduces several custom tokens.

We need to download the IMDB large movie reviews from this site: http://ai.stanford.edu/~amaas/data/sentiment/ (http://ai.stanford.edu/~amaas/data/sentiment/) Direct link : Link (http://ai.stanford.edu/~amaas/data/sentiment/aclImdb_v1.tar.gz) and untar it into the PATH location. We use pathlib which makes directory traveral a breeze.

```
In [2]: BOS = 'xbos'  # beginning-of-sentence tag
        FLD = 'xfld'  # data field tag

        PATH=Path('data/large-movie-reviews-dataset/acl-imdb-v1')
```

# Standardize format

```
In [3]: CLAS_PATH=Path('data/imdb_clas/')
        CLAS_PATH.mkdir(exist_ok=True)

        LM_PATH=Path('data/imdb_lm/')
        LM_PATH.mkdir(exist_ok=True)
```

The imdb dataset has 3 classes. positive, negative and unsupervised(sentiment is unknown). There are 75k training reviews(12.5k pos, 12.5k neg, 50k unsup) There are 25k validation reviews(12.5k pos, 12.5k neg & no unsup)

Refer to the README file in the imdb corpus for further information about the dataset.