

# Tokenization (40%)

## Problem ID: mt2p2

Skrifið forrit sem biður notandann að slá inn nafn á skrá sem inniheldur runu af orðum í sérhverri línu. Forritið tilreiðir textann, þ.e. brýtur runurnar upp í tóka með því að nota hvít bil sem afmarkara (e. delimiter). Ef komma, punktur, upphrópunarmerki eða spurningamerki (, . ! ?) kemur fyrir í lok orðs þá eru þessi greinarmerki meðhöndluð sem sérstakir tókar. Gera má ráð fyrir að þetta séu einu mögulegu greinarmerkin og ef þau koma fyrir í inntaksskrá þá er það í lok orðs (ekki hvítt bil er á undan greinarmerki).

**Ekki** er leyfilegt að nota import setningu í lausninni nema `import typing`.

Fyrsta dæmi um inntaksskrá, `data1.txt`:

```
One Two Three
Four Five Six
Seven Eight Nine
```

Annað dæmi um inntaksskrá, `data2.txt`:

```
One Two Three? Four
Five, Six Seven Eight!
Nine Ten.
```

Þriðja dæmi um inntaksskrá, `data3.txt`:

```
This is a text,
with some
punctuation characters! attached
to some words.
Right?
```

## Inntak

Inntakið er nafn skrár sem á að tilreiða, t.d. `data1.txt` eða hvaða annað skráarnafn sem er sem endar á `.txt`. Inntaksskrá inniheldur 1 til 10 línur, þar sem sérhver lína inniheldur runu af 1 til 15 orðum, og lengd sérhvers orðs er 1 til 10 stafir. Hér er orð skilgreint sem runa af hástöfum og/eða lágstöfum, mögulega með kommu, punkti, upphrópunarmerki eða spurningarmerki í lok orðsins.

Athugið: Forritið ykkar á ekki að skoða sérstaklega hvort þessi inntaksskilyrði standist, né hafna öðru inntaki. Þetta er bara ykkur til upplýsinga um það inntak sem er prófað í prófunartilvikunum. Þið þurfið sem sagt ekki að búast við inntaki utan þessara takmarkana.

## Úttak

Ef ekki er hægt að opna inntaksskrána þá er ekkert úttak myndað. Annars samanstendur úttakið af eftirfarandi:

1. Fjöldi orða í sér línu og sérhvert orð úr inntaksskránni í sér línu þar á eftir.
2. Fjöldi tóka í sér línu og sérhver tóki úr inntaksskránni í sér línu þar á eftir.

**Sample Input 1**

data1.txt

**Sample Output 1**

9  
One  
Two  
Three  
Four  
Five  
Six  
Seven  
Eight  
Nine  
9  
One  
Two  
Three  
Four  
Five  
Six  
Seven  
Eight  
Nine

**Sample Input 2**

data2.txt

**Sample Output 2**

10  
One  
Two  
Three?  
Four  
Five,  
Six  
Seven  
Eight!  
Nine  
Ten.  
14  
One  
Two  
Three  
?  
Four  
Five  
,  
Six  
Seven  
Eight  
!  
Nine  
Ten  
.

**Sample Input 3**

data3.txt

**Sample Output 3**

13  
This  
is  
a  
text,  
with  
some  
punctuation  
characters!  
attached  
to  
some  
words.  
Right?  
17  
This  
is  
a  
text  
,  
with  
some  
punctuation  
characters  
!  
attached  
to  
some  
words  
.  
Right  
?