

<1>

how to identify new page?

· regex

<5>

· check first char of each line
and last char

<4>

check for <> = " "

<aa>

check if middle is a digit

<1>

strip all punctuation ⊕ white space

↓
old spaces before and after each word

<1> hello my name is ~~is~~ jonathan

key word

<1> hello w my w name w is w jonathan w

· jonathan

check if string contains key word w/ space

w jonathan w