→ <u>Word Tokenization</u> :          Words ≈ Token

→ ⓘ have a can opener, but ⓘ can't open the cans.

  ↳ Word Token : 11 words

  ↳ Word Type : Different realization of word.

              10 word Types -

→ <u>Issues w/ W. Tokenization</u> :

  → Finland's → Finland, Finlands, Finland's, Finlands'

  → What've, I'm, shouldn't → What have

                            I am, Should not

  → <u>San Francisco</u>

  → m.p.h → ✓

→ <u>Handling Hyphenation</u> :

  → End-of-line   Hyphen : show-time, initia-time

  → <u>Lexical Hyphen</u> :

      Prefix : eg : co-, pre-, meta-, multi-

→ <u>Sententially Determined Hyphenation</u> :

  → Case-based, hand-delivered

→ three-to-five-year plan.

→ <u>Normalization</u> : U.S.A & USA

→ <u>Case-folding</u> :

    ↳ Reduce all text to lower case.

    ↳ US (us)

      ( ↳ among us there is an imposter
      ( ↳ the us is a financial superpower.

    ↳ Exceptions (Task Dependent)

      ↳ ML , GM

→ <u>Morphology</u> : (Morphemes)

        <u>Uncondiontionally</u>

→ <u>Lemmatization</u>: Reduces inflections or variant forms
               to base forms.

    → car, cars, car's, cars' → car

→ Morphemes are divided in 2 categories :

  → <u>Stem</u> : Core meaning bearly unit.

  → <u>Affixes</u> :

→ **Affixes**:

   → Prefix : un-, anti-, Sanskrit (u-, ati-, pra-)

   → Suffix : -ity, -ation

   → Infix : Abso-damn-lutely

            Un-freakin-believable

→ **Stemming** : Crude chopping of affixes.
             It is always language dependent.

→ **Porter's Algorithm** :