

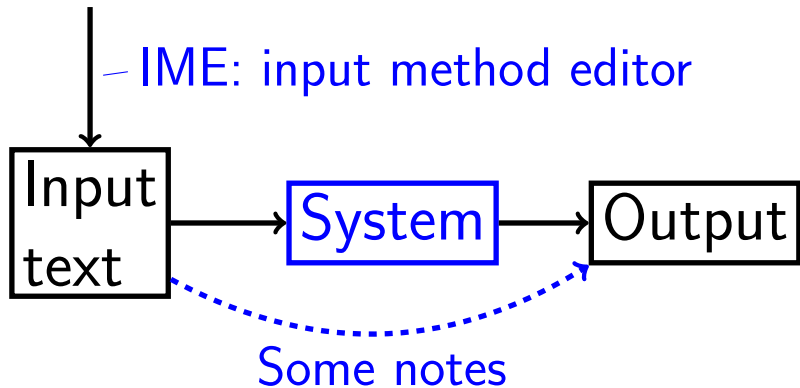
# Some notes on Japanese T<sub>E</sub>Xt Processing

KUROKI Yusuke

`kuroky(at)users.sourceforge.jp`

October 24, 2013

# Overview



# IME: input method editor

- ▶ There are several ways to input Japanese into computer. Usually,
  1. input *kana* first (directly, by romanization, by pocket bell style, by flick input<sup>1</sup>, etc.), then
  2. change them to *kanji-kana-majiri* correctly by human
- ▶ The software, IME, helps both operations above
- ▶ Users freely to choose where they change *kanas* to *kanji-kana-majiri*.
- ▶ Users often turn on IME to input Japanese & off to Latin. In writing T<sub>E</sub>X source, we change the modes frequently.

---

<sup>1</sup>With help of Moe Masuko

# $\text{T}_{\text{E}}\text{X}$ -related systems to operate Japanese

- ▶ De facto standard in Japan:  
     $\text{pT}_{\text{E}}\text{X}$  (engine extention) + jsclasses class files
- ▶ New age:  $\text{LuaT}_{\text{E}}\text{X}$ -ja (macros of  $\text{T}_{\text{E}}\text{X}$  & Lua for  $\text{LuaT}_{\text{E}}\text{X}$ )
- ▶ Experimental stage?:  $\text{ConT}_{\text{E}}\text{Xt MkIV}$
  
- ▶  $\text{upT}_{\text{E}}\text{X}$   
    (change the internal operations of  $\text{pT}_{\text{E}}\text{X}$  into Unicode)
- ▶  $\text{ConT}_{\text{E}}\text{Xt MkII}$  +  $\text{pT}_{\text{E}}\text{X}$
- ▶ CJK package + Takayuki YATO's package
- ▶  $\text{X}_{\text{J}}\text{T}_{\text{E}}\text{X}$ + Takayuki YATO's package

## Note for line-breaks

- ▶ Roughly speaking, Japanese words could be split anywhere due to line-ending
- ▶ Input (e.g., in case of 5 em line-breaking):

これは僕が 飼っている 犬です。	v.s.	This is the dog which I keep.
------------------------	------	-------------------------------------

- ▶ Output:

No Good これは僕が飼っている犬です。

**Good** これは僕が飼っている犬です。

v.s. This is the dog which I keep.

- ▶ Sometimes, we need a little space as the author indicates, e.g., pTeX は中野 賢さんほかにより作られた。

## Note for Unicode input

When we use JIS X 0208 character set, we could sort out which areas are for Japanese and which for Latin easily.

- ▶ multi-byte area should be for Japanese
- ▶ ASCII area should be for Latin

§    § (input \S before Unicode age)  
“    “ ( ‘ ‘ )  
”    ” ( ’ ’ )

In Unicode age, since some signs and marks are combined, we will need indicate which area is in which language.