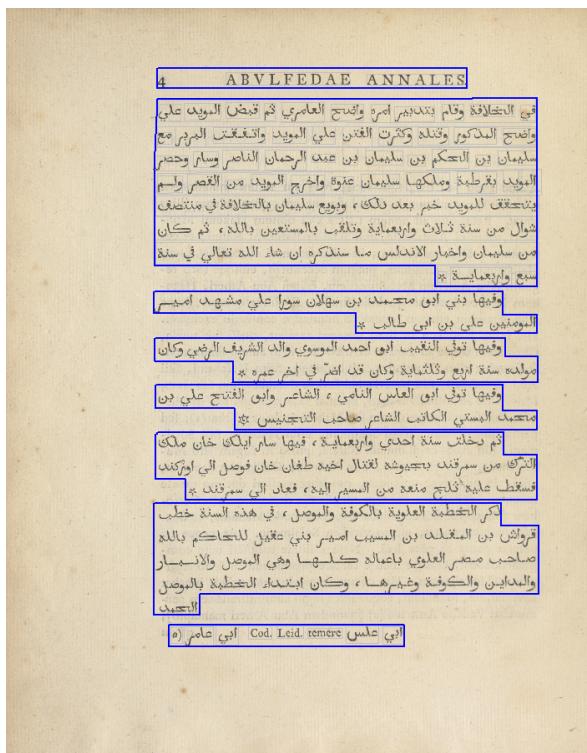


## ALTO Languages support

### 1. Arabic sample (with latin secondary language)

<http://gallica.bnf.fr/ark:/12148/bpt6k62219939/f14.image>

1789-1794, monography



### Language and script

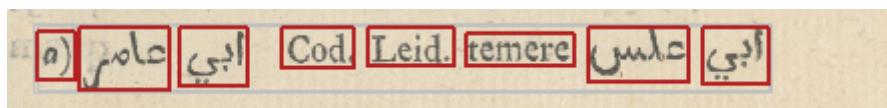
Should be defined at the top level (ALTO root) or at element (block, line, word) level.

Language: arab

Script: naskh

Notes:

1. Most of the time, setting the language at the top level is enough. But some languages use mixed layout modes (Chinese, Japanese), and some documents are multilingual. Eg: this foot note has arab and latin words.



2. Language and script are 2 different things. Eg. Turkish language can be written with latin or arabic scripts.

Use cases: OCR correction, automatic language processing, accessibility, etc.

## Block flow direction

Block flow direction is defined by the logical blocks reading order. This can be implemented in ALTO with:

1. order of XML elements within the file,
2. IDNEXT attribute set on Block element,
3. a top level attribute could also document the block flow direction. Here, top-to-bottom (horizontal lines of text).

Note : for use and reuse of ALTO contents, 1. and 2. are easier to handle than 3. alone (3. implies a topological analysis of the ALTO elements).



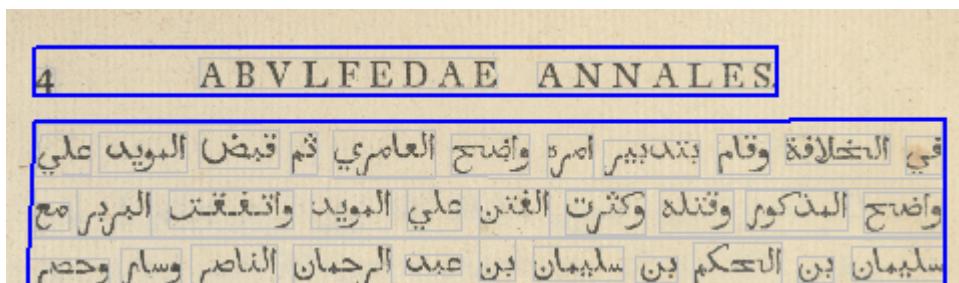
## Reading direction

Reading direction within a line of text can be known :

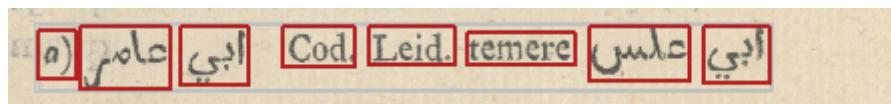
1. with XML elements order within the ALTO file
2. with an IDNEXT attribute on String
3. with a dedicated attribute (ltr, rtl) on Line element
4. with the language information (but there are exceptions).

Use cases and software may have various heuristics to handle this, but 1. and 2. will always be easier to implement than (3) or a language based assumption system (4).

In the sample, first block: LANG="latin" (ltr). Second block: LANG="arab" (rtl)

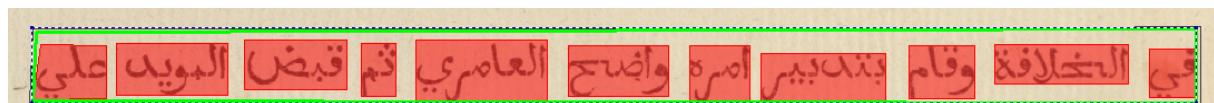


In some cases, reading direction can't be guessed from language. The foot note below, in a mainly arabic page, has a western reading direction (ltr).



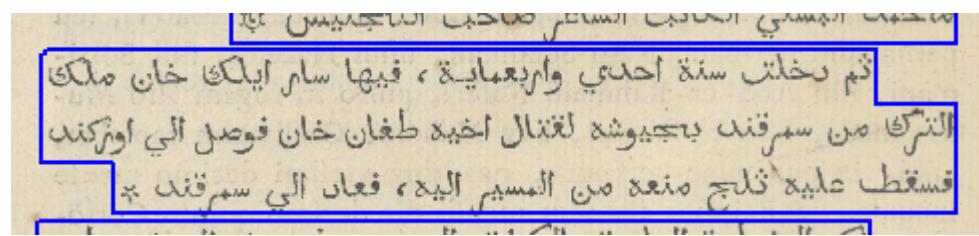
ALTO Use cases:

- **OCR edition/correction** tool with a “Next” button to edit words in a line (or lines in a block).  
This functionality needs the “reading direction” information.



- **Ebook conversion:** the reading direction information is needed to correctly typeset the text.  
First indent of the paragraph, text alignment, etc.

Arabic sample:



Text rendition for arab: right alignment, first indent       

xxx xxx xxxx xxxx xxxx xxxx xxxx TSRIF       

xx xxx xxxx xxxx xxxx xxxx xxxx xx xxxx xxxx xxxx xx

TSAL xxxx xxxx xxxx xx xxxx xxxx xxx

Default rendition (western writing mode) without reading direction information:

      xxx xxx xxxx xxxx xxxx xxxx TSRIF

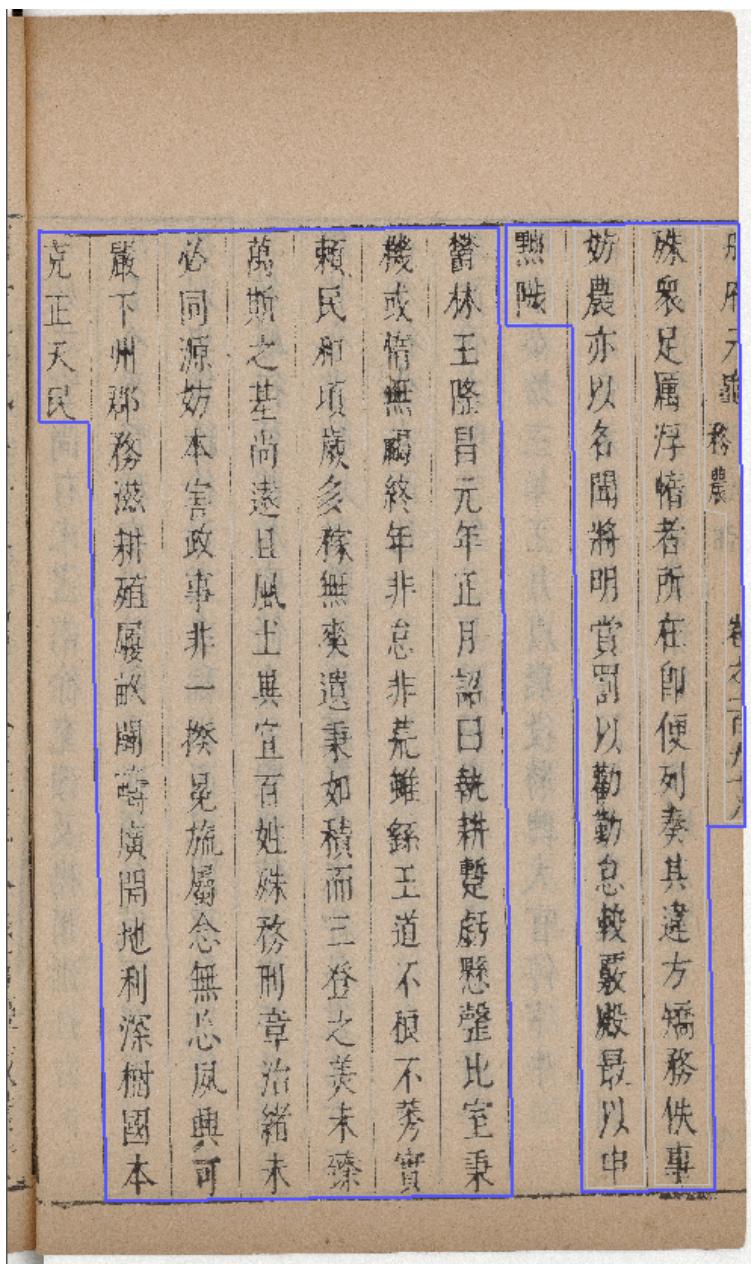
xxx xxxx xxxx xxxx xxxx xxxx xx xxxx xxxx xxx xx

TSAL xxxx xxxx xxxx xx xxxx xxxx xxx

- **Speech synthesis:** order of words in a line must be known

## 2. Chinese samples

### 1. Manuscript



#### Language and script

Could be defined at the top level (ALTO root) level (it's a monolingual document).

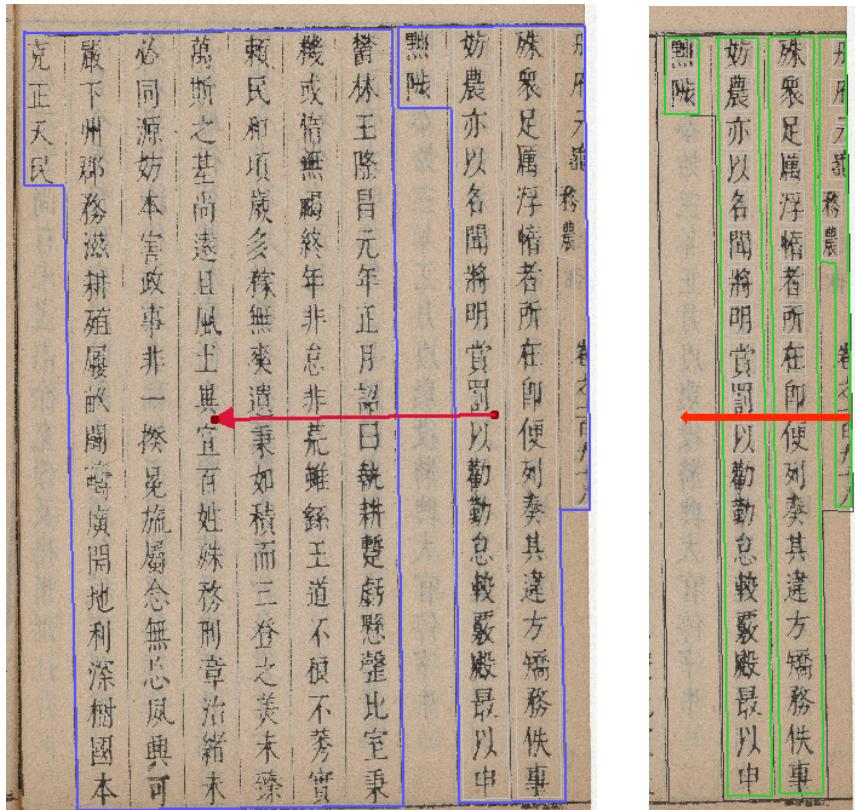
Language: Chinese

Script: Traditional chinese

#### Block flow direction

Block flow direction is defined by the blocks reading order. A top level attribute could also say: (vertical)right-to-left.

This information also rules the way software must process lines of text within a block. Same as above: IDNEXT attribute on Line element, XML elements order, language based guessing.



## Reading direction

Here, reading direction may be set on the top level (the document has an uniform layout).

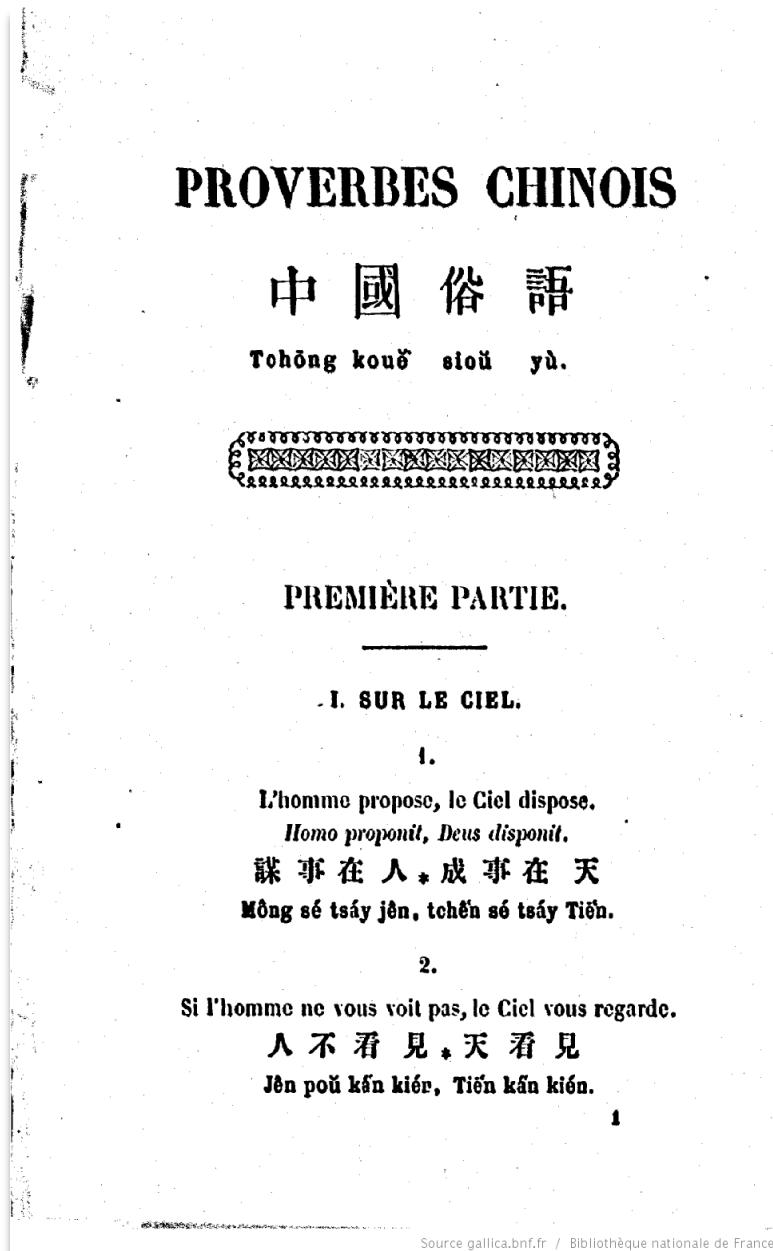
As Chinese can have 2 principal layouts, reading direction must be set to “top to bottom”.



### *ALTO Use cases: same as Arabic sample*

## 2. Mixed languages monography

<http://gallica.bnf.fr/ark:/12148/bpt6k57497215/f12.image>



Source gallica.bnf.fr / Bibliothèque nationale de France

### Language and script

Western layout with French and Latin texts, Chinese text and its phonetic transcription (pinyin):

- Top level : LANG="fr", reading direction: ltr
- Latin texts: LANG="lat"
- Chinese ideograms lines : LANG="chi"
- Phonetic transcription (pinyin): LANG="pinyin"

### Reading direction

- Top level : reading direction= ltr
- or on Chinese texts only reading direction= ltr

Here, reading direction can't be deduced from language (Chinese can be left-to-right or top-to-bottom).

## Other layout features

### Orientation of text

Specify the orientation of text within a (vertical) line.



ALTO already has a Rotation attribute defined on blocks, but use cases are different: a content has a different orientation (relatively to the principal orientation in page or document):

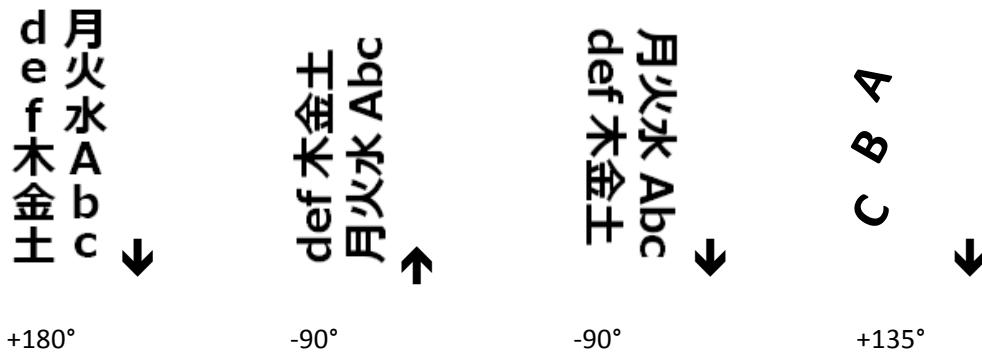
- a large table or an illustration with a sideways layout:



- text within a block with a specific orientation :



Text orientation defines the direction of glyphs relatively to the reading direction (clockwise angles):

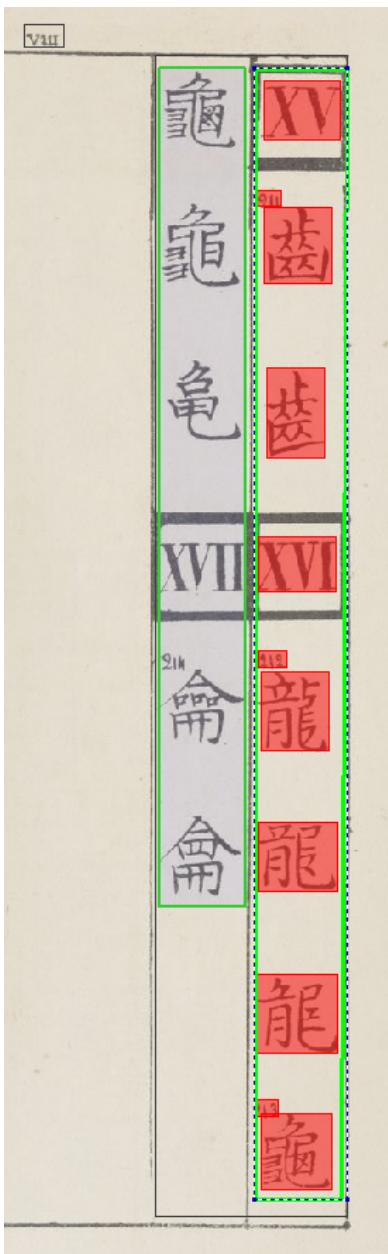


*ALTO Use cases:* text orientation is needed to produce accurate facsimiles.

### Text combine

Monography with combination of multiple characters into the space of a single character. In East Asian documents, the ‘text-combine: horizontal’ effect is often used to display Latin-based strings.

<http://gallica.bnf.fr/ark:/12148/bpt6k6580409q/f14.image>



ALTO Use cases: not clear...

### 3. Conclusion

For core OCR purposes, ALTO functionalities are enough: blocks order, LANG attribute, UNICODE characters.

For advanced use cases, ALTO lacks a complete writing modes model. Data needed to feed this model can be set:

- Before OCR (to help it) – Eg: this block is detected (by human or image analysis system) as horizontal in a Chinese newspaper (mainly vertical).
- During OCR – Eg: this block has been tagged as latin.
- After OCR – Eg: enrichments from other processes are input for “enhanced” ALTO: text, text and layout attributes. Example: ALTO → EPUB → ALTO

### ALTO proposals

- Language and Script:
  - attributes at all levels: page, block, line, word (optional attributes)
- Block flow direction:
  - attribute at the page level (optional): (horizontal)top-to-bottom, (vertical)right-to-left, (vertical)left-to-right, mixed (eg for Chinese and Japanese newspapers).  
NB: bottom-to-top writing mode doesn't exist (in the real world)
  - IDNEXT at block and line levels
- Reading direction:
  - attribute at block and line levels: top-to-bottom, bottom-to-top, right-to-left, left-to-right
  - IDNEXT at word level
- Text orientation:
  - Specific new attribute (angle values) at line and word levels

+ guidelines for languages support with ALTO