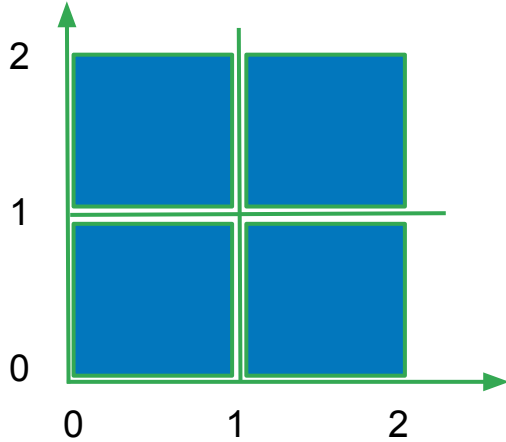Google

# 2. Architecture and Data Structures

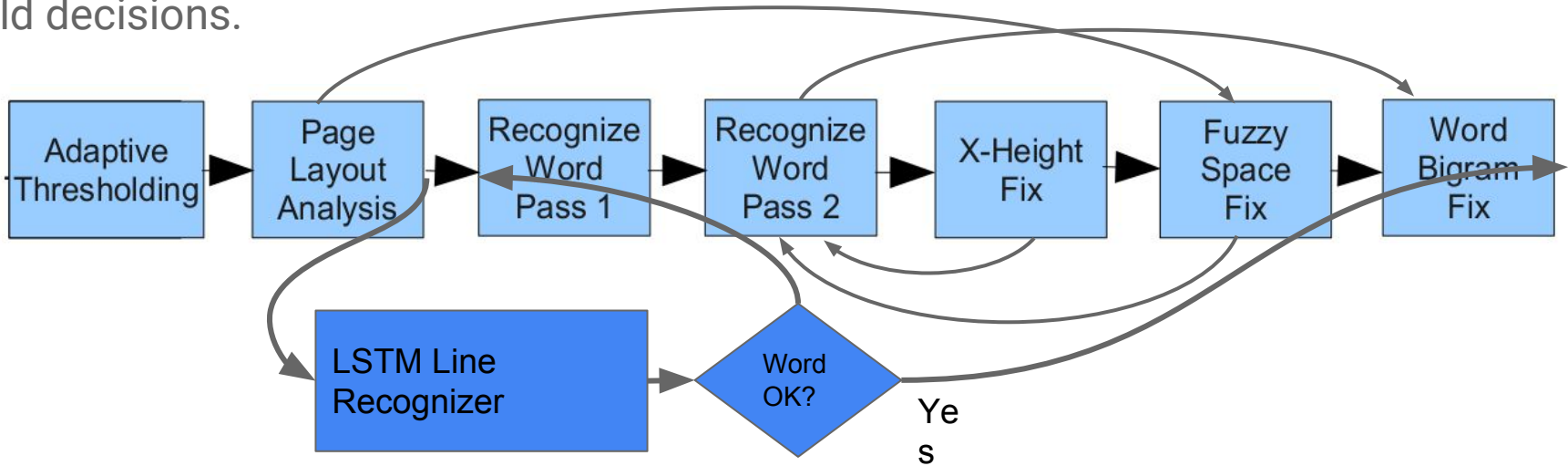A quick tour of the Tesseract Code

*Ray Smith, Google Inc.*

# A Note about the Coordinate System

- The pixel edges are aligned with integer coordinates.
- (0, 0) is at **bottom-left.**
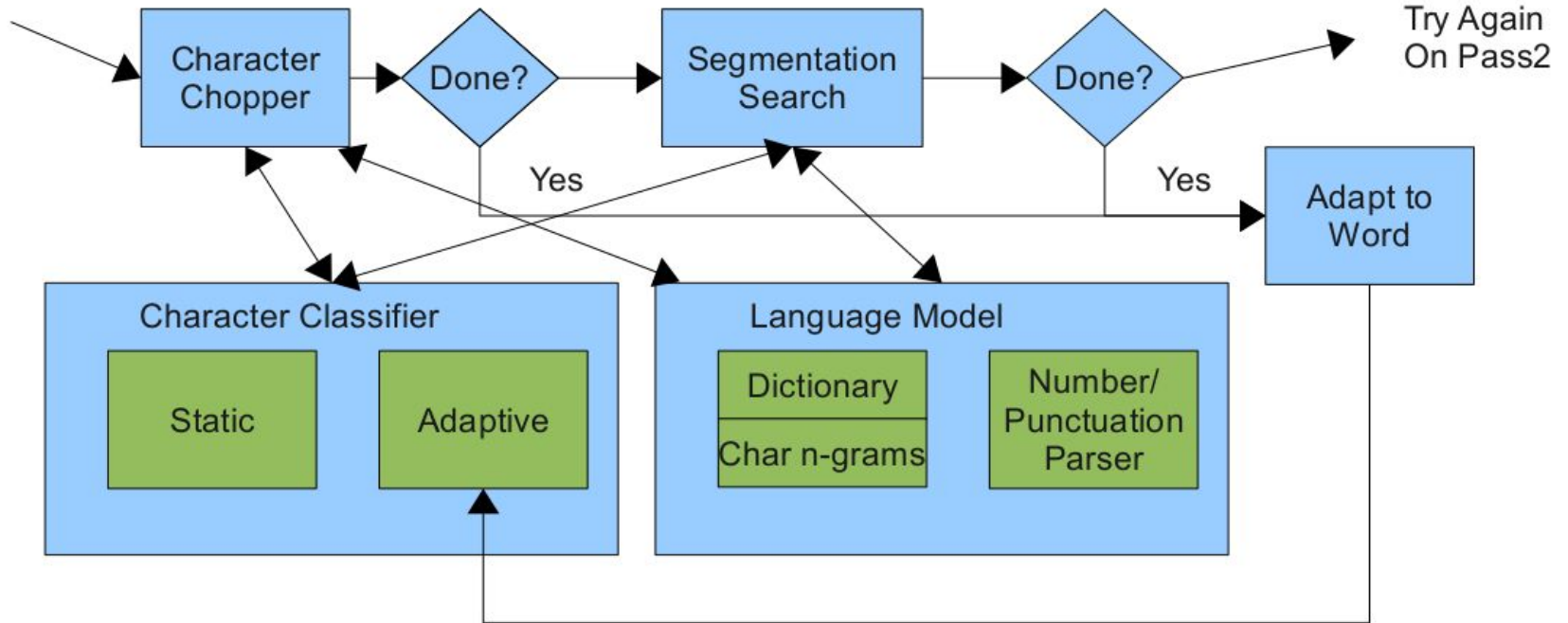- Width = right - left => no silly +1/-1.

Google

# Tesseract System Architecture

Nominally a pipeline, but not really, as there is a lot of re-visiting of
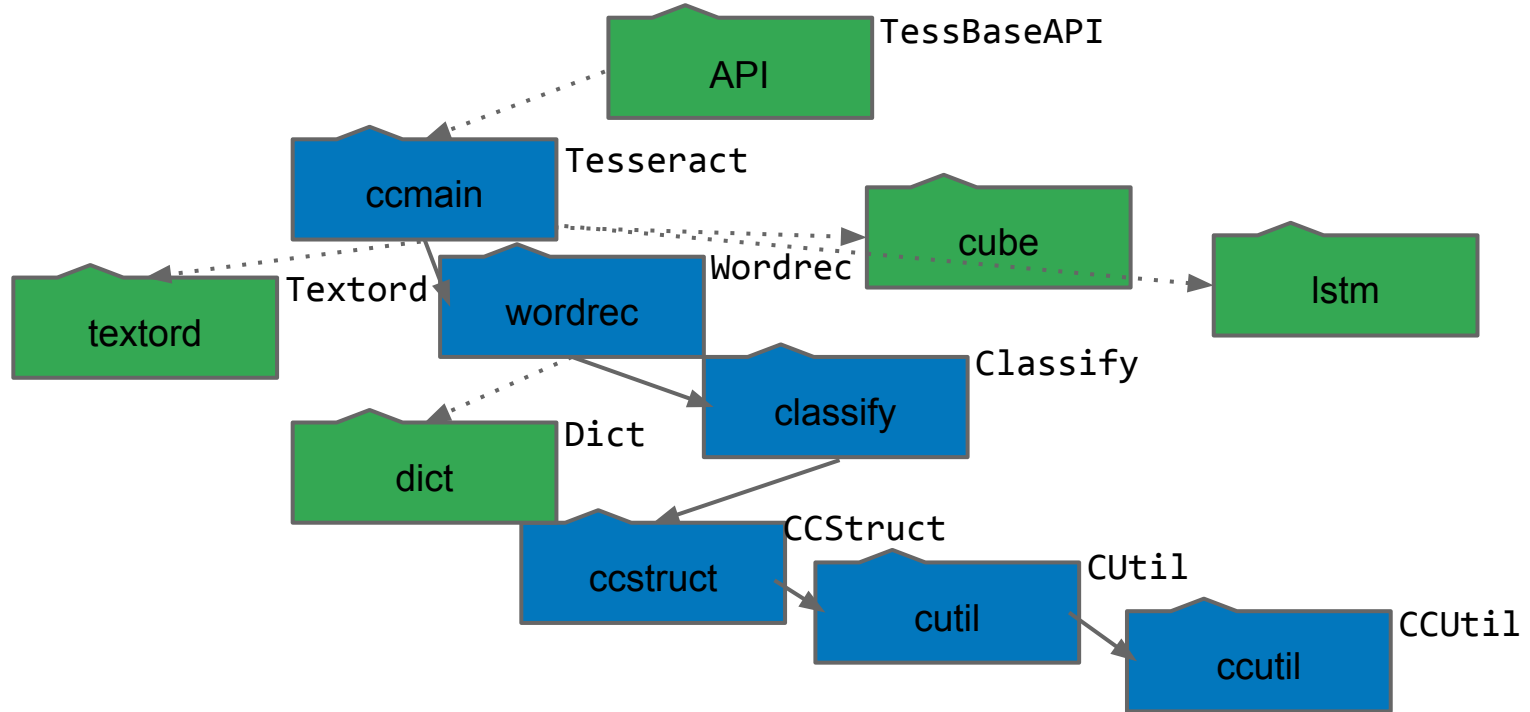
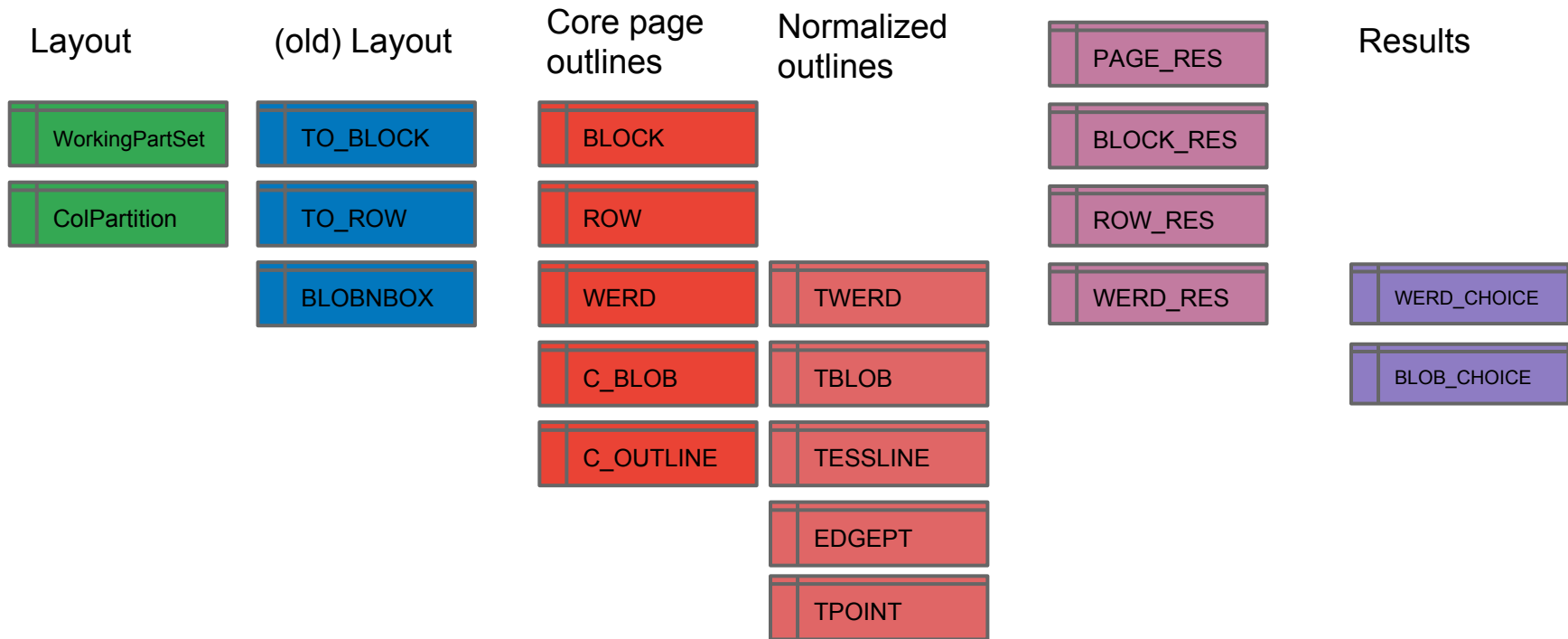old decisions.

# Tesseract Word Recognizer

# The 'C' Legacy

- Large chunks of the code written originally in C.
- Major rewrite in ~1991 with new C++ code.
- C->C++ migration gradual over time since.
- Majority of global functions now live in a convenience directory structure class. (For thread compatibility purposes.)
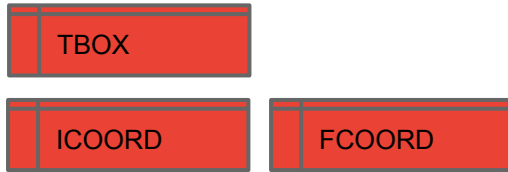
Google

# Directory Structure ~ Functional Architecture

# Key Data Structures = Page Hierarchy

**Layout**

- WorkingPartSet
- ColPartition

**(old) Layout**

- TO_BLOCK
- TO_ROW
- BLOBNBOX

**Core page outlines**

- BLOCK
- ROW
- WERD
- C_BLOB
- C_OUTLINE

**Normalized outlines**

- TWERD
- TBLOB
- TESSLINE
- EDGEPT
- TPOINT

- PAGE_RES
- BLOCK_RES
- ROW_RES
- WERD_RES

**Results**

- WERD_CHOICE
- BLOB_CHOICE

# Software Engineering - Building Blocks

**Coordinates**

TBOX

ICOORD

FCOORD

**Containers**
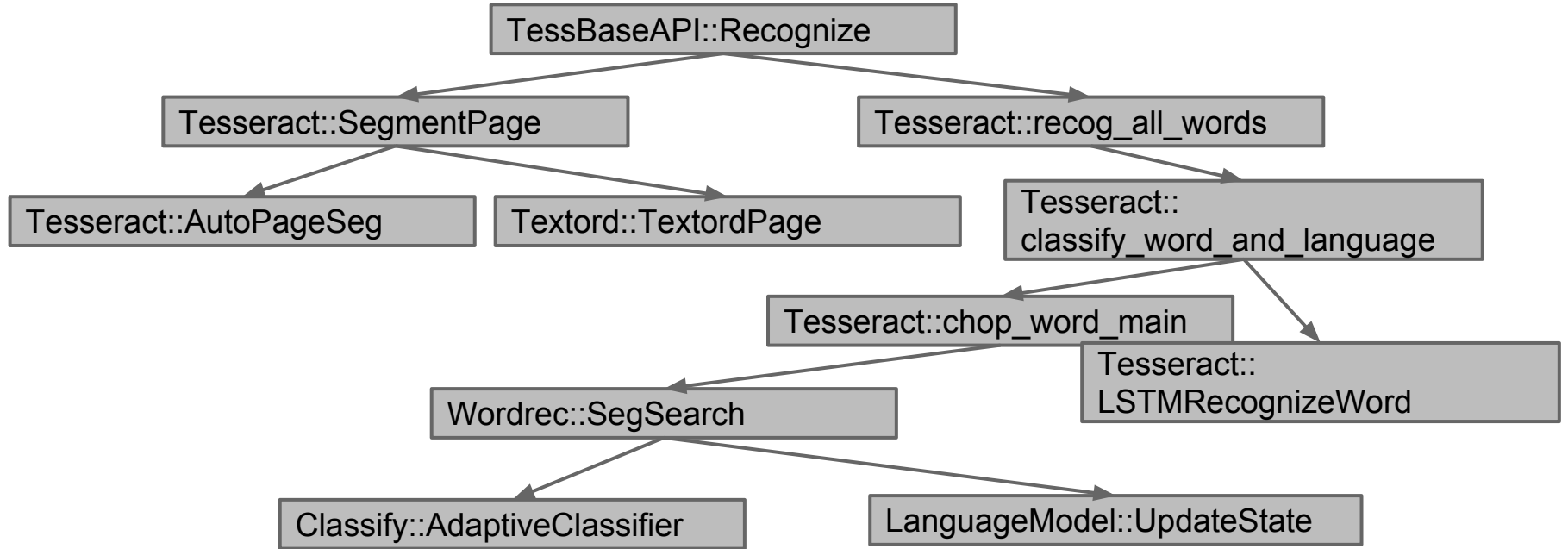
GenericVector

ELIST

CLIST

**Text**

STRING

UNICHARSET

Google

# Key Parts of the Call Hierarchy

# Tesseract's List Implementation

- Predates STL
- Allows control over ownership of list elements
- Uses macros instead of templates

# List Example

Google

# TessBaseAPI : Simple example

Main API class provides initialization, image input, text/hOCR/PDF output:

```
TessBaseAPI api;
api.Init(NULL, "eng");
Pix* pix = pixRead("phototest.tif");
api.SetImage(pix);
char* text = api.GetUTF8Text();
printf("%s\n", text);
delete [] text;
pixDestroy(&pix);
```

Google

# TessBaseAPI : Multipage example

```cpp
TessBaseAPI api;
api.Init(NULL, "eng");
tesseract::TessResultRenderer* renderer =
  new tesseract::TessPDFRenderer(api.GetDatapath());
api.ProcessPages(filename, NULL, 0, renderer);
const char* data;
inT32 data_len;
if (renderer->GetOutput(&data, &data_len)) {
  fwrite(data, 1, data_len, fout);
  fclose(fout);
}
```

# ResultIterator for getting the real details

```cpp
ResultIterator* it = api.GetIterator();
do {
  int left, top, right, bottom;
  if (it->BoundingBox(RIL_WORD, &left, &top, &right, &bottom)) {
    char* text = it->GetUTF8Text(RIL_WORD);
    printf("%s %d %d %d %d\n", text, left, top, right, bottom);
    delete [] text;
  }
} while (it->Next(RIL_WORD));
delete it;
```

Google

# Thanks for Listening!

# Questions?