

Knowledge in Frame: KiF

Natural Language Processing research usually focuses on the rule engine, either through machine learning techniques, which infer rules from already parsed sentences, or through symbolic approaches, where rules are directly designed by linguists. However, when dealing with actual texts, some of the problems that arise during grammar development are often out of the boundaries of the linguistic formalism. Rules are not designed, for instance, to record information or to count occurrences of words or categories. This part of the processing is usually implemented with external scripts, often in Perl or in any other script languages that come as handy for the task, with as a result, a system which might lack both in portability and in maintainability. Furthermore, these external scripts are implemented to serve specific needs, which might not prove easy to extend.

Xerox Incremental Parser (i.e. XIP)

For years, linguists have tried to implement syntactic parsers through complex and constrained theoretical models, such as LFG, GPSG or HPSG ([Gazdar 85]), which despite their efforts ended up in almost intractable systems. Most of these earlier parsers would consider a grammar as a bag of rules, which the rule engine would combine on the fly to analyze a sentence. However, when the grammar comprises up to 10,000 rules, the complexity is such that the system either returns no output or thousands of solutions, usually after a pretty long processing time. To remedy this situation, linguists have implemented new approaches which use probabilities to sort out rules before any processing according to weights computed out of annotated corpora (such as the Penn Tree bank [Grinberg 94]). However, if machine learning parsers have become prevalent today in the research community, they present some problems in an industrial environment. First, the learning stage induces a very specific bias due to the initial corpus on which the parser was trained, which usually diminishes the overall performance on new corpora. Second, if the parser commits the same recurrent errors, the only way to correct them is to present the system with as many examples as possible in order to bias the system out. Finally, the annotation process, whose role is often ignored or minored in the community, introduces many discrepancies and contradictions, which the machine learning model must take into account when extracting and weighting down the rules.

The Xerox Incremental Parser (XIP [Aït-Mokhtar 2002]) on the other hand is still based on a symbolic approach. Grammars are implemented by hand, with the help of a Java IDE, which provides all that is necessary to debug a grammar. Rules in a XIP grammar are grouped in layers, which are applied in sequence, the output of a layer being the input to the next. Thanks to this strategy, the system does not drift into intractable analyses and more importantly always yields a unique solution. XIP also exposes different sorts of rules to solve different linguistics problems, which again is a difference with previous models, which would only provide linguists with one single formalism, that would have to be tweaked to implement rare or exotic structures. Also, compared to probabilistic models, the correction of a grammar is much simpler as it consists of modifying or adding new rules. Grammars for eight languages (French, English, German, Italian, Spanish, Portuguese, Dutch and Japanese) have been developed so far, usually by one linguist, with an average of six months to implement a new language from scratch. In

comparison, annotating a large corpus by hand requires usually a quite large team of students (the task is quite tedious to say the least) over a period of one or two years.

XIP Processing

A syntactic analysis includes the following stages:

1. Tokenization, which splits a document into a collection of tokens along spaces or punctuations.
2. Morpho-analysis, which interpret each token as a word or an expression and returns their part of speech and a set of features.
3. Tagging, which reduces the part of speech ambiguity of a given word.
4. Chunking, which groups together words around a syntactic head
5. Dependency extraction, which extracts the relation between the heads in a sentence.

Example:

If we analyze the sentence: *The bird flies*.

Then each stage will perform the following operations:

1. Tokenization: [the,bird,flies]
2. Morpho-analysis: [the+determiner,bird+Noun, {flies+Noun,flies+Verb}] *flies is ambiguous*
3. Tagging: [the+determiner,bird+Noun, flies+Verb] *flies is no longer ambiguous*
4. Chunking: NP:[the bird] VP:[flies]
5. Dependencies: SUBJECT(flies,bird) and DETERMINER(bird,the).

For each of these stages, XIP provides specific implementations. For instance, tokenization and morpho-analysis are done with transducers, while tagging, chunking and dependency extraction are executed with hand written rules.

Formalism Enrichment

But rules are blind and deaf: they are executed in a “launch and forget” mode, which makes it quite difficult for the linguist to write rules that would benefit from previous analyses, as [Declerck 2002] shows in his article. In order to solve this issue, we first introduced variables (strings, floats, vector and maps) into the formalism to track information throughout the whole linguistic process. We also added the possibility to implement small functions or procedures, which could be called from any rules.

Python Integration

However, if the enrichment of this formalism would eventually turned out to be too complex from within the parser code, the way functions were implemented proved to be a very interesting starting point. XIP is implemented in C++ and function calls were implemented as a specific class, which could be easily derived. Instead of trying to amend again and again the *XIP scripting language* code that was too closely knotted into the rule engine, we moved to another solution, which was to build an interface between XIP and Python.

Python is written in C and provides a straightforward API to create basic objects such as strings, integers or maps. We linked Python and XIP in two ways ([Roux 2006]). First, we transformed XIP into a Python library, which could be loaded from within a Python program. Second, we implemented a specific interface so that Python functions could be executed from within rules. We used the comment section of Python (a set of lines between two triple quotes) as a placeholder for XIP function declarations, which were bound with the actual Python functions implemented in the same file. Thus, every call to *a function* in XIP resulted into a call to the eponymous function in Python. Parameters are translated on the fly from XIP into Python, while the *return* data from Python functions are transformed back into XIP values. A Python function can then be used to return a numerical value that could be exploited from within the rules.

However, even if this solution proved quite reliable, we faced some issues. First, Python cannot have direct access to the internal linguistic structures in XIP, and can only apply to duplications of these structures through a set of dedicated functions. For instance, a syntactic node is known to Python only as an integer *id*, which is used to fetch its features through a specific function. Worse, memory leaks proved an absolute nightmare to track down.

Second, since XIP communicates with Python through a C programmatic interface, the XIP binary is heavily dependent on the version of Python with which it is linked. If XIP is linked with Python 2.6.2, then running grammars with Python functions embedded requires the exact same Python version on another machine or Python binary libraries such as *math* will not load at all.

Knowledge in Frame Language

Despite the issues mentioned above, Python had proven very useful in many cases, and since we wanted to keep a script language in our linguistic arsenal, we just could not simply put Python aside without replacing it with something else. We decided to implement a new language that would be close to Python in some aspects, but would be more suitable to our linguistic purposes. The Knowledge in Frame language (KiF) was implemented with the objectives of having a full access to all XIP linguistic structures with a formalism that would be as light and simple as Python. Thus, KiF can modify XIP variables or linguistic structures and KiF functions can be transparently accessed from any rules without an interface, as with Python. Furthermore XIP can be used to run KiF scripts, which allows for a pre-processing of documents before passing the result to the rule engine itself. Actually, the integration is so transparent, that KiF functions can be implemented anywhere in a XIP grammar.

Informal Description

KiF is a multithreaded language implemented in C++, and it is available on most platforms (Windows, Linux or Mac OS). It is a cross-over of C++, Python and Java. From Python, we have kept the notion of all-purpose dictionaries and vectors, while from C++ and Java we have borrowed the notion of type declaration. As Java or Python, it is a garbage-collector based language, where each object has a specific reference, which indicates whether the object has been fully released or not. We don't use the Python indentation style but rather the more traditional C++ or Java's utilization of curly brackets.

Example

```
int i;      //we declare an integer
i+=10;     //we increment this integer by 10

//we loop
for (i=0;i<10;i++) {
    int j=i*2;
    println(i,j);
}
```

Contextual Evaluation

From the above example, one would wonder what would differentiate KiF from Java or C++. The most important difference resides in the way expressions are analyzed. In KiF, expressions are evaluated from the left to the right, and the *type* of the expression is given by the leftmost element. For instance, strings, when used in a numerical environment are translated into the number they represent. Thus, the string "10" is automatically translated into the number 10. In the same way, any integer or float is also automatically translated into a string if the context demands it.

```
int i;
string s="10";
i=7+s;      //the result is: 17
s+=i;       //the result is: "1017"
```

All types have different interpretations according to their context. A vector in a numerical expression returns its length. A file as a string, returns a pathname, while as an integer, it returns its size on the disk.

Parentheses

Another big difference is the way parentheses are interpreted. In most languages, parentheses are used to isolate or gather specific elements in a numerical expression: $(10+i)*5$. In KiF, they create a local context, where the new interpretation is imposed by the first element of the expression. $(10+i)*5$ is still interpreted as in other languages, since 10 is a number. However, the interpretation of an expression such as: "MY "+ $(10+2)$ is quite different. The parenthetic part $(10+2)$ is no longer interpreted as a string operation but as an integer operation, thanks to the 10. Thus, "MY "+ $(10+2)$ is: "MY 12", while "MY "+ $10+2$ is: "MY 102".

Below is a list of expressions, whose interpretation is quite specific to KiF:

```
int i="2"+"3"+5;      // i= 2+3+5=10
i=("2"+"3")+5;        //i= ("23")+5=28
string s="MY "+(10+2); //s="MY 12"
s="MY "+10+2;         //s="MY 102"
```

File Evaluation

The small program below shows another example of this contextual evaluation. If we loop in a file with an integer, then KiF will read the file character by character. If we loop in the same file with a string, then each iteration will yield a full line.

```
file f('myfile.txt', "r");
int i;
for (i in f)
    println(i, i.chr()); //the file is read byte by byte, we display the equivalent character
...
string s;
for (s in f)           //the file is read line by line
    println(s);
```

The operator *in* is re-interpreted according to its first argument.

Marshalling, un-marshalling

The main reason for this choice is that in NLP we manipulate strings at length. Actually, NLP could be described as the science of string processing. For example, the vector $[1, 2, 3]$ can be very easily transformed into the string `"[1,2,3]"`, but if conversely a string is fed to a vector, then KiF will assume that this string is a vector description and will parse it to re-generate the corresponding object. This *marshalling* and *un-marshalling* applies to most objects in KiF. Complex structures can then be stored in databases as strings, leaving the task of parsing the objects back to KiF. It becomes much simpler to store a lexicon in a database. For instance, the feature list, which is rather tedious to describe in tables, can then be simply recorded as a string, which will be automatically re-analyzed as a map, when needed. The same idea presides over the choice of translating a string into a number when in the context of a calculus.

```
//s is initialized with a map description, it becomes the string: '{"gender":"masc","pers":2}'
string s='{gender:"masc","pers":2};
map m=s; //m is initialized with s, it becomes the map: {'gender':'masc','pers':2}
```

Frames

Frames are in KiF the equivalent of classes in other languages. Their declaration is very similar to a class definition in Java or in C++.

```
frame test {
    int i;
    string s;
    function _initial(int k) { //the constructor is always _initial
        i=k;
    }
}
//To create instances of that class is pretty straightforward:
test toto(10);
//And access is also quite simple
println(toto.i);
```

Frames can be sub-framed and methods or operators can be overloaded. It is also possible to call the super-method of a mother class if necessary.

Dynamic Cast

Frame objects are also interpreted in context. We have described above how a simple string can be transformed into an integer, if this string finds itself in the context of a calculus. Frame objects can also be subjected to the same transformations. KiF provides a similar mechanism as the *toString* method in Java, requiring developers to implement specific functions to yield the right interpretation. In KiF, these functions bear the same name as the required cast. Thus, a developer can implement a *string()*, a *int()* or a *float()* function at this effect.

```
frame test {
    int i;

    //A simple integer cast function
    function int() {
        return(i);
    }
}

test t;
t.i=10;
int k=20+t; //k is 30...
```

This function will be automatically called by the interpreter in the context of a calculus. If no function is available, KiF will return *null*.

In the same way, it is also possible to implement a *frame* converter to translate a given object into another frame object. Consider for instance the following example, which implements the frame *othertest*.

```
frame othertest {
    int k;
    //A cast into a test object
    function test() {
        test t;           //we create our test object
        t.i=k;            //we copy the value, in which we are interested
        return(t);        //we return this new object
    }...
}

othertest o;           //we create our object
o.k=100;               //we initialize o.k
test t=o;              //now t.i=100, the converter is automatically called to transform an othertest into a test.
```

Associate Functions (binding)

Associate functions are another specificity of KiF. Any variable can be bound with a specific function, which is called every time this variable is modified. This associate function must have two parameters (or three if it is a frame), which are basically the current value of the variable and its new value.

//before is the current value of the variable, and after its new value

```
function modify(int before, int after) {println(before,after);}
```

```
int o with modify;  
o=10;    //will trigger a call to modify...
```

Multithreading

Multi-threading is usually considered, with some good reasons, to be difficult. Since KiF is a language, which is targeted to linguists with little knowledge about programming, we wanted to keep multithreading as simple as possible. Actually, in KiF to declare a function as a thread consists of replacing the keyword *function* with *thread*. The whole language is implemented in such a way that access to variables or call to other functions is automatically protected within any *threads*.

```
//a simple thread description...  
thread mythread(int k) {  
    print(k, ",");  
}
```

Synchronization

However, *threads* must synchronize with each other sometimes. KiF provides a full arsenal of tools at this effect: *mutex*, *join*, *semaphore*, and *synchronized variables*, yet implemented in such a way, that their utilization does not require any complex programming. For instance, the function *wait* puts a thread in hold, while the method *cast* releases it. The *wait* function accepts as parameters, a simple list of strings: *wait("one", "two", "three")*. A *cast* on one of these strings will release all the *threads that wait on this string*. The *wait* function will then return as a result the very string that was cast in the first place.

```
//a thread waiting on strings...  
thread mythread(int k) {  
    string s=wait("test", "off");  
    println(s+", "+k); // print: "off,10"  
}  
  
mythread (10);    //We launch it  
pause(0.001);    //We wait for 10ms, in order for the thread to execute the wait  
cast("off");      //We release it
```

Synchronized Variables

Synchronized variables are used in conjunction with the function: *waitonfalse*, which has been implemented as a semaphore. Any type of variables can be used as a synchronized variable, as long as it has a Boolean interpretation (*it can be a frame object*). The variable must be associated with the pre-defined function: *synchronous*.

```
//a synchronized variable  
int sync with synchronous=3;  
...  
waitonfalse(sync);
```

Thus, each time *sync* is modified, the system tests its Boolean value. When *sync* is 0, it then releases *waitonfalse*.

Join Threads

However, using *synchronized variables* is often too complex. A more simple way of synchronizing a set of threads with their master is to declare these threads as *join*, and then wait on their completion with *waitonjoin()*.

```
join thread toto(int i) ...
```

```
toto(10); //we launch our threads
toto(20);
waitonjoin(); //and wait for their completion...
```

Furthermore, the *waitonjoin* works at the thread level, which means that many *waitonjoin* can run in parallel waiting for different sets of threads in as many master threads.

XIP integration

The other important part of this language is the way it communicates with XIP. The bulk of this adaptation is a list of specific KiF types which are mapped over linguistic objects. For instance, syntactic nodes are mapped over *node* objects, which can then be used to access the part of speech, the lemma or any features from a given node. It is also possible to modify the content of that node, since we are dealing with an encapsulation of the node internal pointer in KiF.

Example:

For instance, the following XIP rule, which computes a passive subject between two words (each word is a syntactic node denoted as a *#d* variable), also calls a KiF function to test some constraints on the *verb* node.

XIP side

This rule applies on a sentence such as “*The cat is chased by the dog*” with the following dependencies computed so far:
Complement(chases,dog)
PrepObj(dog,by)

Our XIP rule tests whether a verb has a complement introduced with a “by” preposition. It then checks if the verb is in the passive mode.

if (Complement(Verb#1,Noun#2) & Prepobj(Noun#2,Prep#3[lemma:“by”]) & *testpassive*(Verb#1)) *subjectpassive*(#1,#2).

KiF side

```
function testpassive(node n) {
  map m=n.data(); //we load the features from the node n
  if (m["mode"]=="passive") //if the mode is passive, we return true
    return(true);
  return(false);
} //The return value will then be used in the XIP rule to either activate or deactivate its execution
```


Of course, the above example could have been implemented with the adequate XIP formalism on features, but it illustrates the simplicity with which XIP linguistic information can be transparently passed onto KiF. All the linguistic objects have their particular mapping in KiF, with all their specific data available for testing or modification. Thanks to KiF, specific treatments, which are algorithmic by nature, can be naturally intertwined into the fabric of a grammar.

Evaluation

We have compared KiF with Python in terms of speed and memory footprint. We did two evaluations: one as a standalone language, the other one from within XIP.

As a standalone interpreter

We have implemented many Python programs over the years, most of them doing intensive string manipulations in order to clean and extract proper data for our linguistic tools. We re-implemented about 20 of them into KiF. From our test, we could show that Python has a memory footprint which is about 15% to 20% smaller than KiF, with a speed which is between 10% and 15% faster.

Within XIP

We have also re-implemented many of our previous Python procedures into KiF, with results quite different from the previous tests. Python is still about 5% less greedy than KiF, in terms of memory footprint. As for the speed, the difference is quite difficult to compute, since most of these procedures account for only a tiny fraction of the whole processing time. Still, from our tests, it appears that Python and KiF are quite similar in terms of speed, thanks to the reduced data duplication, which consumes both time and resources.

KiF vs. Python

From our experiments, KiF appears slightly slower and greedier than Python. Actually, the difference is not that large that any KiF scripts would ruin the overall rule engine efficiency. KiF does not represent a terrible drift in terms of memory and speed, when evaluated from within XIP. However, we believe, that the benefits of improvements achieved in terms of functionality outweigh the differences in speed and memory footprint, since KiF represents a much better language in terms of integration into XIP than Python, as all linguistics objects and structures are seamlessly accessible from any scripts. Also, KiF resolves some of the problems that plague script development with dynamic languages. First, KiF variables must be declared. Non-declared variables, means that errors such as misspelled names, can only be detected at run time, while in KiF they can be detected at compile time. Second, variable types are fixed, which means that variable types do not evolve during run time, reducing the risk of dark corners as in Python where problems can jump out of rarely used sections of code. We still think that Python is a remarkable language, with very powerful features, but the lack of declarativity makes it difficult to debug and to maintain. Furthermore, Python programs cannot be debugged in development environments such as eclipse when they are run from within XIP. In the case of KiF, we have implemented an internal debugger, which can be called from within the rule engine at any times. Linguists can stop the parsing and examine which structures and which variables were sent to their KiF scripts, to see how these structures and variables are actually processed.

Conclusion

Natural language processing is the art of string manipulations. A large part of the effort of implementing grammars or NLP-based applications requires the programming of complex string transformations to deal with a large variety of inputs, either as documents or as linguistic data. The implementation of a language such as KiF, which provides a natural marshalling-un-marshalling embedded into the fabric of the language helps simplify these tasks a lot. The language is of course richer than the simple description that has been given in this article. For instance, it is possible to build external libraries in which new types can be encoded to enrich the language. We have already encapsulated external libraries such as SQLite or FLTK into KiF modules, which can be loaded from any programs or grammars.

The language still needs some polishing and some improvement both on speed and on memory footprint, in the future we plan on improving KiF's speed and reducing its memory footprint.

Reference

- Gazdar G., Klein E., Pullum G., Sag A. I., 1985. *Generalized Phrase Structure Grammar*, Blackwell, Cambridge Mass., Harvard University Press.
- Tapanainen P., Järvinen T. 1994. *Syntactic analysis of natural language using linguistic rules and corpus-based patterns*, Proceedings of the 15th conference on Computational linguistics, Kyoto, Japan, pages: 629-634.
- Grinberg D., Lafferty John, Sleator D., 1995. *A robust parsing algorithm for link grammars*, Carnegie Mellon University Computer Science technical report CMU-CS-95-125, also Proceedings of the Fourth International Workshop on Parsing Technologies, Prague, September, 1995.
- Roux C. 1999. *Phrase-Driven Parser*, Proceedings of VEXTALL 99, Venezia, San Servolo, V.I.U. - 22-24.
- Aït-Mokhtar S., Chanod J-P., Roux C., 2002. *Robustness beyond shallowness incremental dependency parsing*, NLE Journal, 2002.
- Declerck T. 2002, *A set of tools for integrating linguistic and non-linguistic information*, Proceedings of SAAKM.
- H. Cunningham, D. Maynard, K. Bontcheva, V. Tablan., 2002. *GATE: A Framework and Graphical Development Environment for Robust NLP Tools and Applications*, Proceedings of the 40th Anniversary Meeting of the Association for Computational Linguistics (ACL'02), Philadelphia, July 2002.
- Blache P., Guénot M-L. 2003. *Flexible Corpus Annotation with Property Grammars*, BulTreeBank Project
- Roux C., 2004. *Une Grammaire XML*, TALN Conference, Fez, Morocco, April, 19-22, 2004.
- Roux C., *Coupling a linguistic formalism and a script language*, Proceedings of the Third Workshop on Constraints and Language Processing, pages 33-40, 2006.

[Python] <http://www.python.org/>