# The Bioscanners Project at Sourceforge

**Detlef Groth (1), Stefan Müller (1), Joachim Selbig (1)**

**(1): Potsdam University, Bioinformatics Group, c/o Max Planck Institute of Molecular Plant Physiology, Am Mühlenberg 1, 14476 Potsdam, Germany**

Although biologists and bioinformaticians are permanently challenged with parsing and analysing output of various biological tools, often utilising different formats, a standard method for solving such tasks has not been found, yet. To overcome the problems of commonly used standalone applications (difficult data integration), and Bio-frameworks (complex programming interface, slow data parsing), we generated the Bioscanners project. The usage of scanner generators in application coding ensures easy programming, little influence of personal programming styles, a small code base, easy maintenance, and very high processing speed.
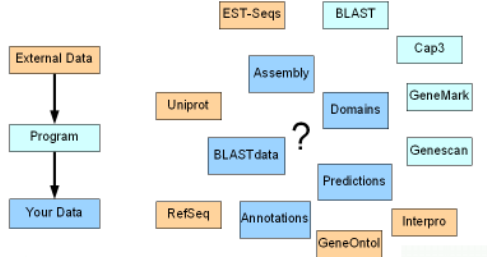


**Figure:** The data integration problem: How to look at all those different formats? How to make relations between all those data?

### Scanner Advantages

- Easier to program than traditional parsers
- Standalone programs
- One file - one application
- No installation, just download and execute
- Almost no library dependencies (just libC)

Besides easy and fast data scanning, the storage of the scanning result in a format, that can be efficiently used in downstream processing and to selectively extract information, is important for the user. For that reason the output of our scanners is standard database code, suitable to be used for SQL compliant databases like PostgreSQL, MySQL or SQLite. That data integration platform ensures the best possible performance in data exploration, especially in a relational context, as well as in maintaining data integrity.
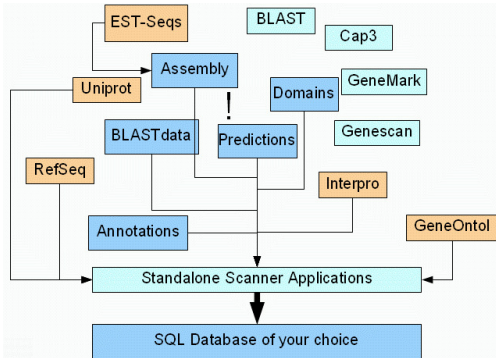


**Figure:** The data integration solution: Write small programs for data scanning with SQL output. That output can be directly piped in a database.

### Database advantages

- All data in one place (even one file with SQLite)
- No API required, just ask your data using SQL
- Easy to learn, within minutes
- Sophisticated GUIs are available for interactive data queries and data exploration

As an example, our BLASTScanner, available from the sourceforge project page, is a single C source file which can be compiled on any modern computer platform to a small 20Kb executable not depending on any external library or runtime. The data processing time required of the scanner is shorter than the time required by the database to import the data. Even input files of several hundred MB can be translated into database code within several seconds.
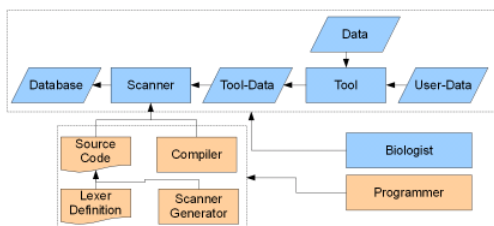


**Figure:** Data-Analysis pipeline using a scanner generator. Programmers and biologists work is strictly separated.

```
$ blastall -p blastn -i sample.fasta -d mydb > sample.blastn
$ BLASTScanner-Linux-x86 --infile sample.blastn \
```

```
    --prefix sam1 | sqlite3 sample.db3
$ sqlite3 sample.db3 "select count(*) from sam1_hits where score > 500"
3428
```

**Figure:** Invocation of BLASTScanner with a redirect to the database application of sqlite3.

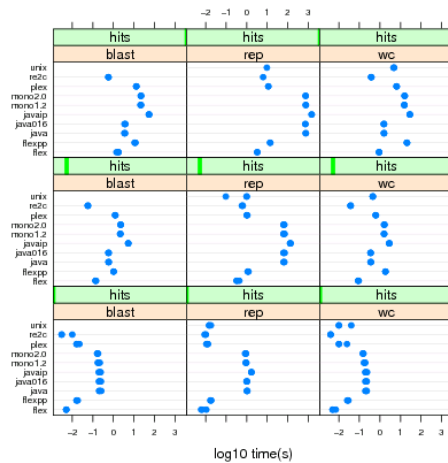| Table 1. Scanners, Compilers and WC binary sizes | | | | | |
|---|---|---|---|---|---|
| We compared different programming languages and scanner generators concering their speed, memory and binary sizes (+++ best, ++ very good, + good, - bad). wc-sizes are size of an word counter implementation in kB. | | | | | |
| scanner | lang. | reference | compiler | wc-sizes | speed/mem |
| Flex | C | Paxson (1995) | GNU gcc 4.1.2 | 16 / 489 | ++ / ++ |
| Flexpp | C++ | Paxson (1995) | GNU g++ 4.1.2 | 21 / 1002 | - / + |
| RE2C | C | Bumbulis and Cowan (1993) | GNU gcc 4.1.2 | 7 / 7 | +++ / ++ |
| TPLex | Pascal | Graef (1998) | Free Pascal fpc 2.0.4 | 109 / 109 | + / +++ |
| JFlex | Java | Klein (2008) | SUN javac 1.6.0 | 7 / NA | + / - |
| GPLEX | C# | Gough (2008) | Mono gmcs 2.0.1 | 13 / 6525 | - / - |



**Figure:** Time usage for different scanners and tasks. wc=word counter, rep=simple character replacer, blast=simple blast scanner. Input were BLAST files with 1 (left), 100 (center) and 1000 (right) hits. The latter file was around 40Mb large.

### Conclusions

Scanner applications can serve as a platform for the development of standardized bioinformatics related tools. The source code for some proof of prinicple implementations is freely available at the sourceforge project page under an open source BSD-license. Programmers, who would like to participate in the development are highly welcome.
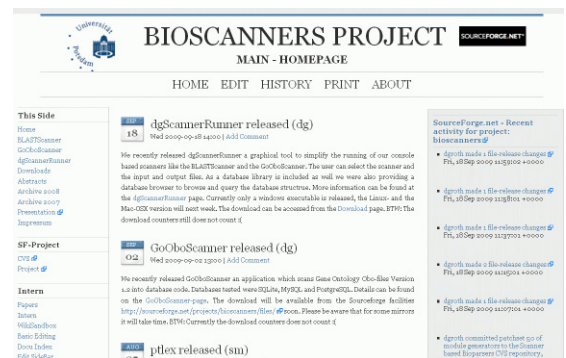


**Figure:** The website of the bioscanners project at http://bioscanners.sf.net

- RE2C generated scanners are single C-code files
- Can be compiled with any modern C-compiler
- No external dependencies, run just standalone
- 20-80kb binaries (Linux, Win32, OSX, Solaris, OSF-1, ...)
- Example: on a normal modern desktop PC, BLASTScanner translates 50MB BLAST-file in one second into database code
- The database actually needs longer to import the code
- Scanners are magnitudes faster than the parser of the Bio*-approaches
- Other scanners for BLAST-M8 files, Gene Ontology Obo-files, gene predictor files are available on the project page as well
- Vizit an join us at: http://bioscanners.sf.net
- **Contact:** dgroth@uni-potsdam.de