# Java DOM Parsers to Convert KGML into SBML and BioPAX Common Exchange Formats

**Kyung-Eun Lee, Myung-Ha Jang, Arang Rhie, Chin Ting Thong, Sanduk Yang and Hyun-Seok Park\***

Department of Computer Science, Ewha Womans University, Seoul 120-750, Korea

## Abstract

Integrating various pathway data collections to create new biological knowledge is a challenge, for which novel computational tools play a key role. For this purpose, we developed the Java-based conversion modules KGML2SBML and KGML2BioPAX to translate KGML (KEGG Markup Language) into a couple of common data exchange formats: SBML (Systems Biology Markup Language) and BioPAX (Biological Pathway Exchange). We hope that our work will be beneficial for other Java developers when they extend their bioinformatics system into SBML- or BioPAX-aware analysis tools. This is part of our ongoing effort to develop an ultimate KEGG-based pathway enrichment analysis system.

*Availability:* You may download the conversion module from Sourceforge.net (https://sourceforge.net/projects/ngssequencealig/). The feedback from developing the conversion tools will provide valuable feedback for the continued evolution of a universal exchange format.

*Keywords:* systems biology, visualization, common exchange format

## Introduction

Recently, various frameworks for the XML formats of biological network information were proposed (Achard *et al.*, 2001; Finney *et al.*, 2006). Among them, two approaches have gained popularity as promising candidates to be adopted as future standards in the field: SBML and BioPAX (Hucka *et al.*, 2003; Strömbäck *et al.*, 2005). On the other hand, we have been working on a KEGG-based pathway analysis system in past years (Jang, 2010; Ham, 2007). Unfortunately, KGML (Kanehi-

*Corresponding author: E-mail neo@ewha.ac.kr
Tel +82-2-3277-2831, Fax +82-2-3277-2306

sa and Goto, 2000) is not a representative common pathway exchange format. The KEGG site (http://www.kegg.org/) provides metabolic pathway data, not only in KGML format but also in BioPAX (http://www.biopax.org/) Level 1 and SBML (http://www.sbml.org/) format. Still, the KEGG site itself does not offer conversion tools. We were able to find only two conversion tools, called KEGG2SBML (http://sbml.org/Software/KEGG2SBML) and KEGGconverter (http://www.grissom.gr/Keggconverter). Both KEGG2SBML (Funahashi, 2004) and KEGGconverter (Moutselos *et al*, 2009) deal only with SBML.

Thus, we decided to develop Java-based KGML2SBML and KGML2BioPAX parsers with a graphical user interface. These parsers not only convert KGML into SBML but also convert KGML into BioPAX.

## Overview of the BioPAX and SBML standards

BioPAX and SBML are becoming one of the most widely used common data exchange formats that can be widely used to simulate and visualize biological networks (Strömbäck *et al.*, 2005). Both BioPAX (Level 3) and SBML (Level 1) can encode signaling pathways, metabolic pathways, and regulatory pathways, although SBML can represent finer details (Fig. 1). However, they
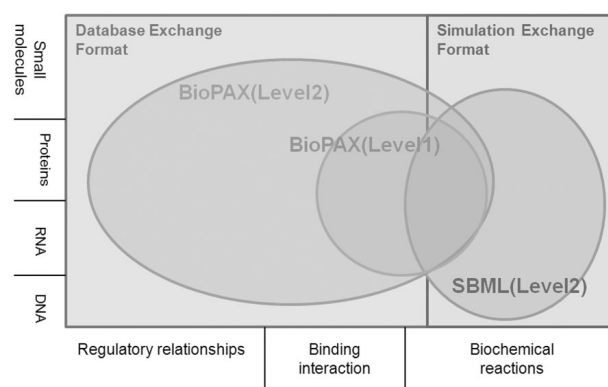


**Fig. 1.** SBML (Systems Biology Markup Language) for simulation model exchange formats, and BioPAX (Biological Pathway Exchange) for database exchange format. Each SBML Level is intended to describe a biochemical network or a pathway. BioPAX is defined in a number of steps, called levels. Level 1 focuses on metabolic networks, and Level 2 adds molecular interaction networks.

are not competing standards. Rather, they are comple-mentary approaches to describe related phenomena.

SBML is targeted at modeling systems of dynamic bi-ochemical reaction networks. BioPAX is an RDF (Re-source Description Framework) and OWL-based stand-ard and assumes a static view of the network model. BioPAX does not contain mathematical formulas but provides more detailed information concerning the in-dividual molecules and interactions. If one's aim is sim-ulation and computational analysis, then SBML has bet-ter API supports than BioPAX. However, BioPAX is a rich standard when it comes to encoding the inter-actions, such as genetics.

Unfortunately, BioPAX and SBML, the two most wide-spread standards for storing pathway data, are semanti-cally different. Conversion modules between KGML and these two formats may lead to loss or distortion of data. This seriously limits the automatic integration of models. To obtain sufficient generality in the model, many details have to be omitted in the model hierarchy.

## Implementation

We implemented the KGML2SBML and KGML2BioPAX parsers, which convert KGML into SBML Level 2, ver-sion 4 documents, and BioPAX. These two parsers are implemented with the Java DOM (Document Object Model) parser and require Java APIs for XML Process-ing (http://java.sun.com/xml/jaxp/index.html). The map-pings between KGML and these two formats are sum-

(a) Mappings between KGML and SBML

(b) Mappings between KGML and BioPAX



**Fig. 2.** The KEGG Pathway element is a root element, and one pathway element is specified for one pathway map in KGML. The entry, Relation, and Reaction elements specify the graph information. SBML is defined in terms of the participating Reactants and Products, along with optional modifier Species, an optional Kinetic Law. The basic building blocks of a BioPAX model are the Physical Entity, Interaction, and Pathway classes. Each Interaction describes relationships between Physical Entities.

(a) The KEGG Markup Language (KGML) is an exchange for-mat of the KEGG graph objects, especially the KEGG pathway maps. The detailed format is described in http://www.genome.jp/kegg/xml/.

(b) A definition of a model in SBML consists of lists of one or more of these various components: *Compartment, Species, Reaction, Parameter, Unit definition,* and *Rule.* The detailed for-mat is described in http://sbml.org/Documents/Specifications.



**Fig. 3.** An exemplary conversion of KGML into SBML.

marized in Fig. 2. KGML2SBML uses the LIGAND database of KEGG as input to generate SBML documents. For example, *KineticLaw* in reaction and specific names of compounds can be added in a converted SBML file with the use of KEGG Ligand information.

An exemplary conversion between KGML and SBML is shown in Fig. 3. Conversion between KGML and BioPAX has been done in a similar fashion. The online SBML Validator (http://sbml.org/Facilities/Validator) has been used to check the syntax and internal consistency for converted SBML documents. Also, BioPAX Validator (http://www.biopax.org/biopax-validator/index.html) has been used to check converted BioPAX documents and fix syntatic and semantic errors. All converted SBML and BioPAX documents are available (https://sourceforge.net/projects/ngssequencealig/).

## Conclusion

Most current simulation and analysis packages already support SBML and BioPAX, and more are in the process of being extended to support them. Also, most of the top journals in the bioinformatics field encourage authors to prepare models of biochemical reaction networks using common exchange formats and to deposit the model into the BioModels database. For this reason, we developed simple Java utility tools for converting KGML to SBML and BioPAX. Making the databases available in common pathway exchange formats would be a first step toward building an open source pathway analysis resource.

### Acknowledgments

## References

Achard, F., Vaysseix, G., and Barillot, E. (2001). XML, bioinformatics and data integration. *Bioinformatics* 17, 115-125.

Finney, A., Hucka, M., Bornstein, B.J., Keating, S.M., Shapiro, B.E., Matthews, J., Kovitz, B.L., Schilstra, M.J., Funahashi, A., Doyle, J.C., and Kitano, H. (2006). Software infrastructure for effective communication and reuse of computational models. Systems modeling in cell biology: from concepts to nuts and bolts. *MIT Press.* pp. 369-378.

Funahashi, A., Jyoraku, A., and Kitano, H. (2004). Converting the KEGG Pathway Database to SBML. *5th Int. Conf. Syst. Biol. (ICSB2004).*

Ham, S.I., Song, E.H., Yang, S.D., Thong, C.T., Rhie, A., Galbadrakh, B., Lee, K.E., Park, H.S., and Lee, S.H. (2007). J2.5dPathway: a 2.5D visualization tool to display selected nodes in biological pathways, in parallel planes. *Genomics Inform.* 7, 171-174.

Hucka, M., Finney, A., Sauro, H.M., Bolouri, H., Doyle, J.C., and Kitano, H. (2003). The systems biology markup language (SBML): a medium for representation and exchange of biochemical network models. *Bioinformatics* 19, 524-531.

Jang, M., Whang, J.W., Lewis, E., and Park, H.S. (2010). An interpretation of biological metabolites and their reactions based on relation degree of compound pairs in KEGG XML Files. *JSW.* 5, 187-194.

Kanehisa, M., and Goto, S. (2000). KEGG: kyoto encyclopedia of genes and genomes. *Nucl. Acids Res.* 28, 27-30.

Moutselos, K., Kanaris, I., Chatziioannou, A., Maglogiannis, I., and Kolisis, F. (2009). KEGG converter: a tool for the in-silico modelling of metabolic networks of the KEGG Pathways database. *BMC Bioinformatics* 10, 324.

Strömbäck, L., and Lambrix, P. (2005). Representations of molecular pathways: an evaluation of SBML, PSI MI and BioPAX. *Bioinformatics* 21, 4401-4407.