

YARN: НАЧАЛО / YARN BEGINS

П. И. Браславский

М. Ю. Мухин

О. Н. Ляшевская

А. А. Бонч-Осмоловская

А. А. Крижановский

П. В. Егоров

Диалог-2013

WordNet для русского языка

- Необходимый ресурс для автоматической обработки текстов (наряду с морфологическим анализатором, большим аннотированным корпусом и т.п.).
- Общедоступный большой ворднет на материале русского языка пока не существует!

Анонсированные ресурсы

- Оригинальные лексические базы:
 - RussNet (рук. И. В. Азарова)
 - Тезаурус РуТез (НИВЦ МГУ).
 - Проекты, связанные с переводом и адаптацией PWN:
 - СПб, Госуниверситет путей сообщения [Сухоногов, Яблонский 2005]
 - Новосибирский ГУ [Гельфенбейн и др.]
- степень разработки?
- незавершенность или закрытость

Пути решения организационных проблем

- Использование существующих источников (словарей, корпусов)
- Автоматические методы
- Crowdsourcing

Yet Another RussNet (YARN)

- Инициативная группа – представители УрФУ, Kontur Labs, ИММ УрО РАН, ВШЭ
- участники из Санкт-Петербурга, Казани, Челябинска, Томска, Петрозаводска.
- http://groups.google.com/group/yarn_org/

Основная идея YARN

- Сочетание традиционных принципов создания ворднетов и wiki-подхода.
- Открытость, открытость и еще раз открытость.

Wiki-принцип

- Большое количество участников (например, студенты-лингвисты).
- Централизованное хранение данных, редактирование тезауруса через веб-интерфейс.
- Контроль процесса редактирования и качества.

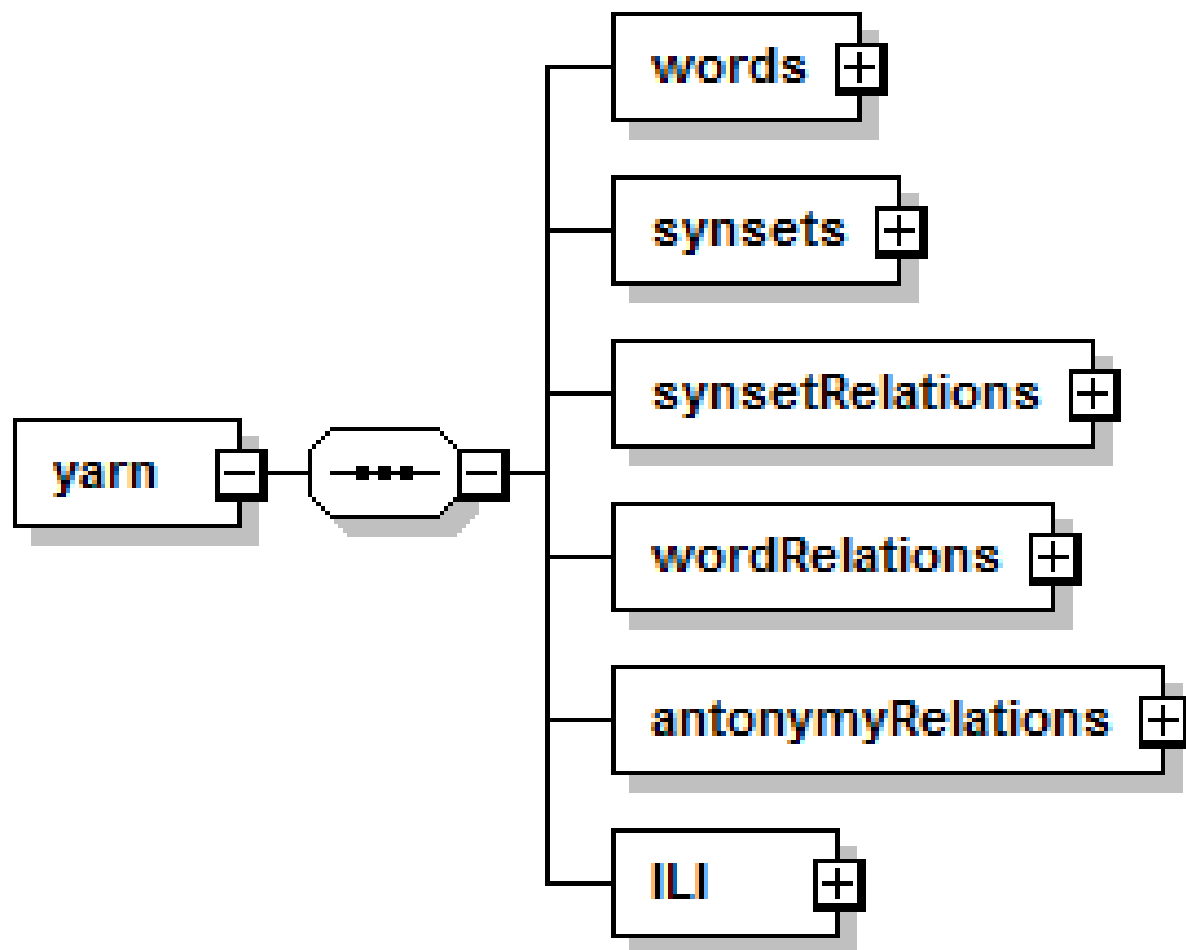
Лингвистические принципы YARN

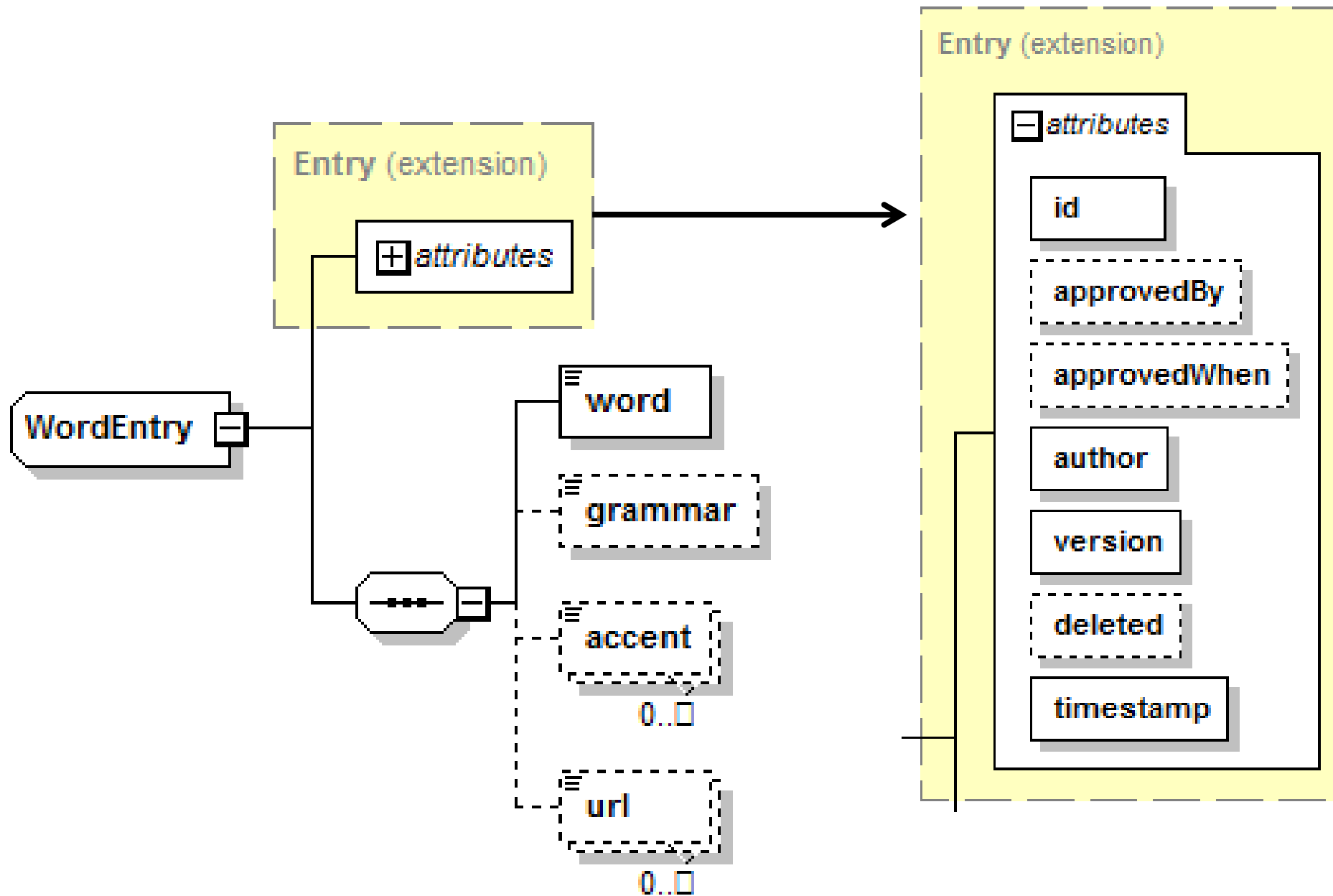
- YARN ориентируется на мировую практику создания wordnet-подобных ресурсов.
- Начальное лексическое наполнение – существительные, прилагательные и глаголы.

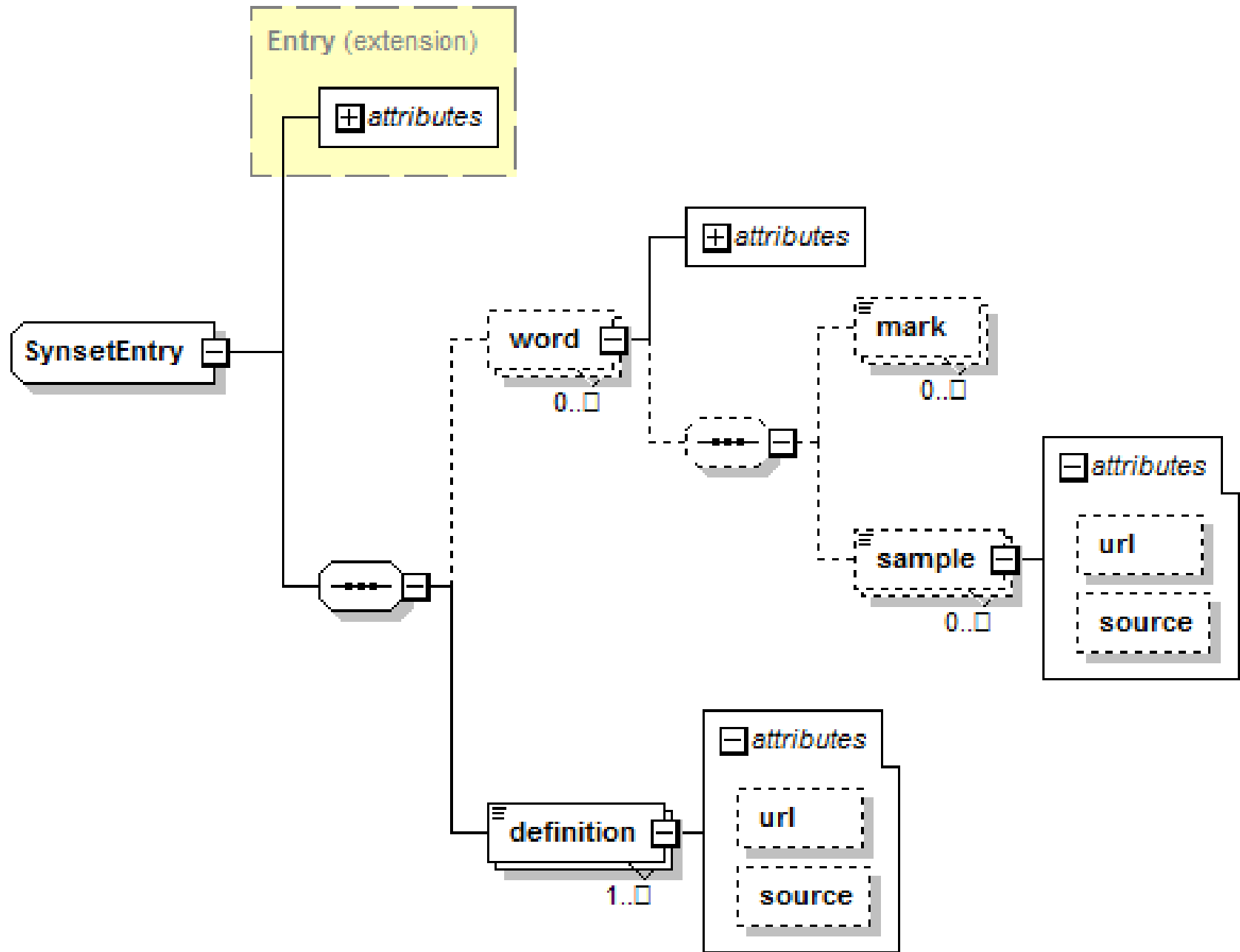
Стандартная WN-подобная организация YARN

- Основа – синсеты (группы синонимов и квазисинонимов).
- Синсеты упорядочены между собой иерархически и связаны отношениями гиперо/гипонимии, антонимии, холонимии/меронимии

Формат XML: компоненты схемы







Первые шаги

- 20 000 частотных существительных (НКРЯ)
- Данные русского Викисловаря
- Сборка синсетов

Приглашаем к сотрудничеству

- <http://russianword.net>
- http://groups.google.com/group/yarn_org/

YARN — проект в стадии становления

- максимальное расширение числа участников редактирования словарных данных
- привлечение возможностей автоматизированного извлечения лексикографической информации
- открытость проекта в плане его дальнейшего развития и использования.

Благодарности

- Исследование осуществляется при финансовой поддержке РГНФ (проект № 13-04-12020 «Новый открытый электронный тезаурус русского языка»).
- Благодарим участников группы *yarn_org* за активность, замечания и предложения.

Спасибо за внимание!

Браславский П.И.

Мухин М.Ю.

Ляшевская О.Н.

Бонч-Осмоловская А.А.

Крижановский А.А.

Егоров П.В.

pb@kontur.ru

mfly@sky.ru

olesar@gmail.com

abonch@gmail.com

andrew.krizhanovsky@gmail.com

pe@kontur.ru

P.S. Princeton Wordnet, started in 1986. <http://wordnet.princeton.edu/>

PRINCETON UNIVERSITY

Search

WordNet

A lexical database for English



What is WordNet?

- What is WordNet?

People

News

Use

Current News

[George A. Miller](#), who began the WordNet project in the mid-1980s, passed away on July 22, 2012 at the age of 92.

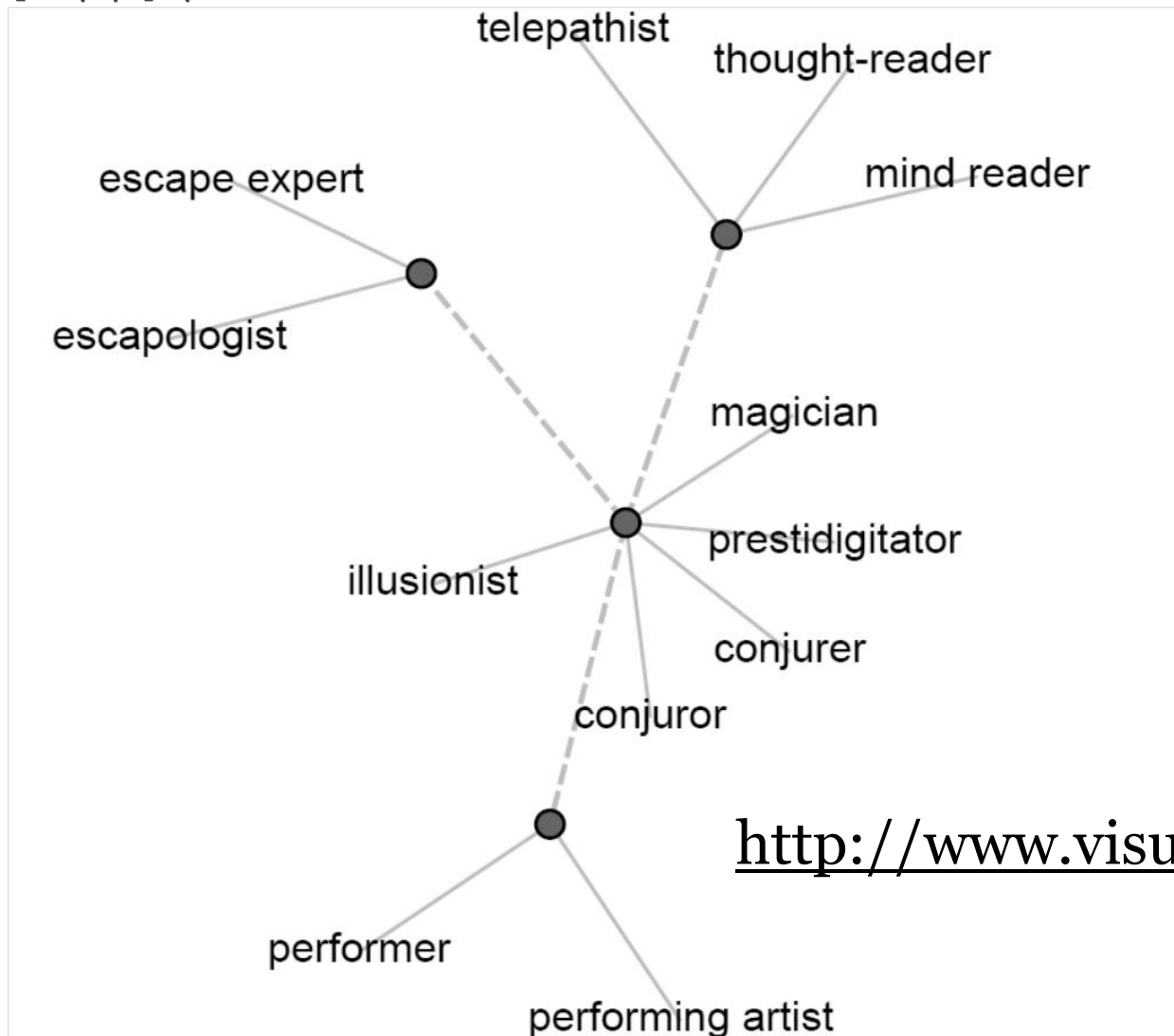
You can read his obituary [here](#).

We appreciate your comments and suggestions, especially when they are constructive and help us improve WordNet. We get numerous questions

Princeton Wordnet (PWN)

- Стандартный подход к организации таких ресурсов.
- Сегодня *ворднетами* называют любые лексические базы, построенные по схожим принципам: они состоят из синсетов (*synset*, от *synonym set*) – «смыслов», которые выражаются набором квазисинонимов. В свою очередь синсеты связаны между собой различными отношениями: гипероним/гипоним, мероним/холоним и др.

Наглядное представление фрагмента PWN



<http://www.visualthesaurus.com/>