

Scrape and graph your way to conference glory: Building the ultimate Call for Papers tool

Lju Lazarevic
Developer Advocate
@ellazal
<http://lju.io>

**This story begins with a
Developer Advocate, facing a
time-honoured challenge...**

***What conference should I submit
to?***

I'm a Developer Advocate at Neo4j



What does a Developer Advocate do?


- Think about problems a developer is trying to solve
- Come up with relevant examples
- Help you with the tools they love!
 - (I really, really love graph databases)

Ways a Developer Advocate might do this:

- Help out in the community
- Write blog posts and create example code
- **Present at conferences and meetups**

Challenges behind finding suitable conferences

- No centralized source
- No obvious way to search for certain topics/tags/keywords
- Sometimes there are no tags/topics/keywords!
- “Conference Driven Development”



This particular Developer Advocate was looking for potential conferences a few months ago...

There's got to be a better way...

Then I had my 💡 moment...

- I did a bit of web scraping many years ago
- I am graph database geek - all about those connections!
- What if I could bring these approaches together, and:
 - Create a Conference Call for Papers tool!



And there's more!

We could use this tool to:

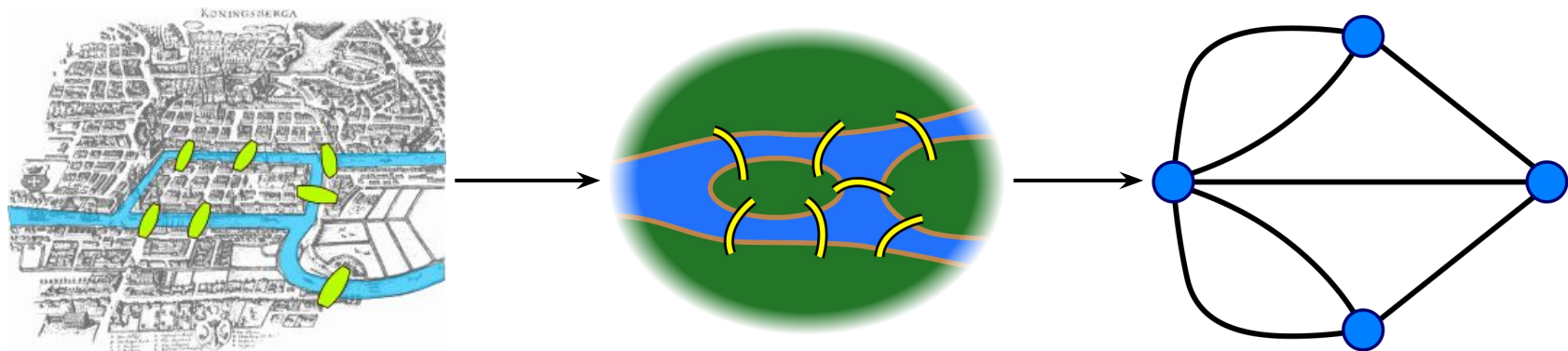
- Get insight into tech trends
- Group similar-themed conferences together
 - Maximise that “Conference Driven Development” effort!
- Build more relevant examples for the developer community
- And so much more!

But first! What is a graph database?



A graph is...

...a set of discrete objects, each of which has some set of relationships with the other objects



Seven Bridges of Königsberg problem. Leonhard Euler, 1735

A graph database is...

...a database that stores data entities and their relationships in a graph structure.

Anatomy of a (property) graph database includes:

Node (Vertex)

- Main element from which graphs are constructed

Relationship (Edge)

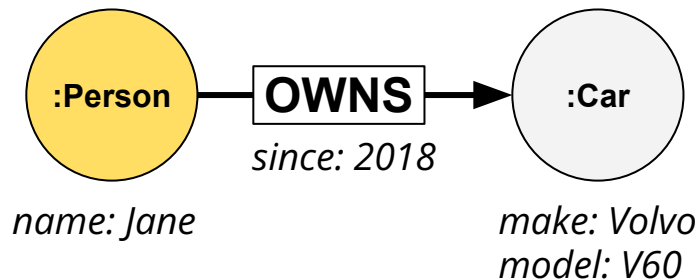
- A link between two nodes, has direction and type

Label

- Define node category

Properties

- Enrich a node or relationship



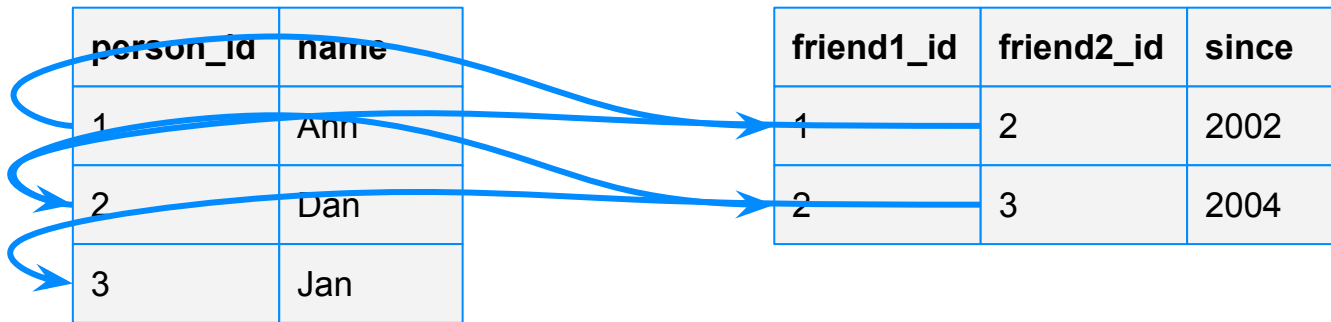
What are the differences between a relational database and a native (property) graph database?



Relational databases - joins on read

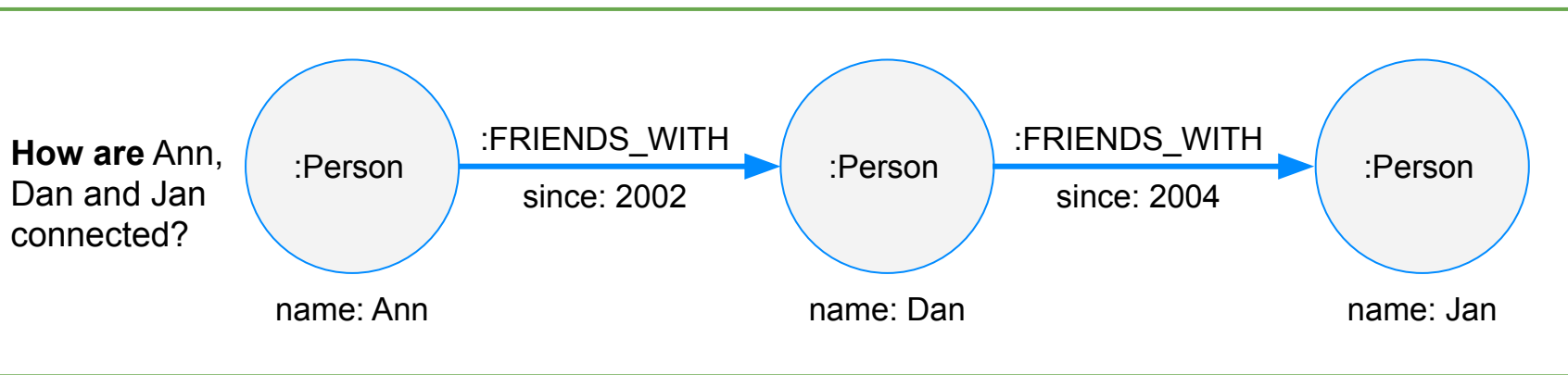
- Typically normalized data & mapping tables
- Joins at query time to reconstitute data and find connections
- Hypothesize on 'is the data connected?'
- Increase in joins → exponential increase in query execution time

Are Ann, Dan and Jan connected?



Native graph database - joins on write

- Relationships are first-class citizens
- Joins at write time - “physically” join connected entities
- Hypothesize on ‘how/why is the data connected?’
- Increase in joins → linear increase in query execution time

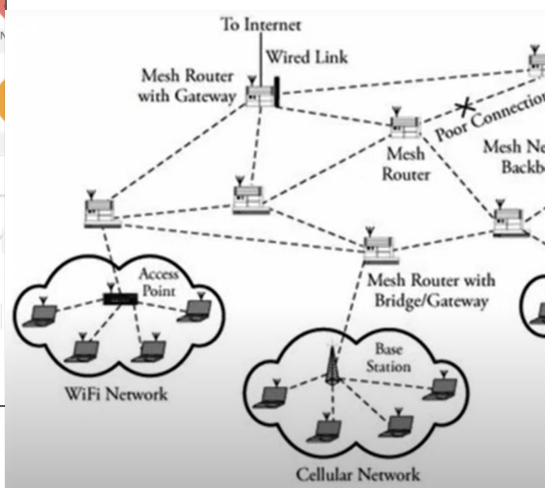


Great graph use-cases

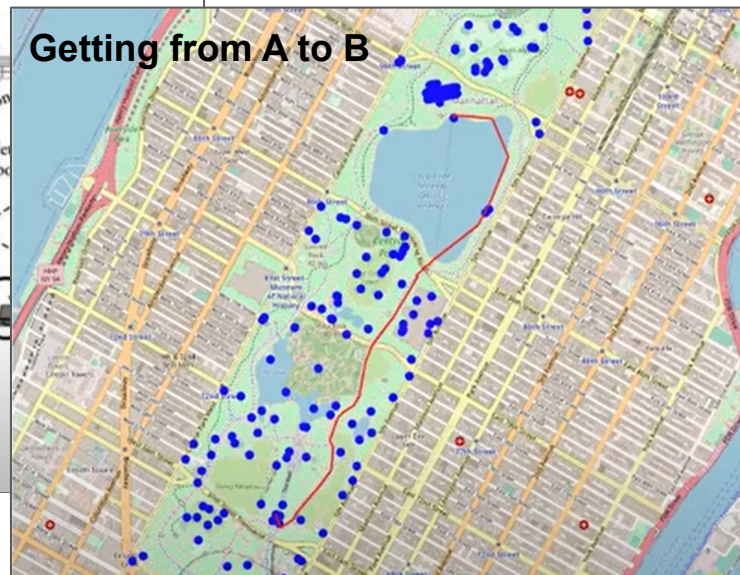
Detecting Fraud



Predicting outages



Getting from A to B



So, how do we build this CfP tool?



Building the CfP tool

1. Identify common places for Call for Papers information
 - a. Web search
 - b. Ask around
 - c. Past experience



The easy way to manage your call for papers


PaperCall enables event organizers to easily manage their call for papers and talk submissions. Speakers are already using PaperCall to manage their events. Are you ready to hear from them?

[Start now >](#)

The smart way to manage Call for Papers, Speakers and Agenda for your event.

Cloud based, safe and easy. Your speakers will love it, too!



 [EVENTS](#) [TALKS](#) [COMMUNITY](#) [PRICING](#) [CREATE EVENT](#) [LOGIN](#) [REGISTER](#)

Events for Tech People

Particip technical events anywhere in the world

[START TODAY](#)

FEATURED EVENTS



 <p>SUN, 05 JUN 2022 9:00 AM</p> <p>Agile Testing Days</p>	 <p>THU, 04 NOV 2021 9:00 AM</p> <p>Code BEAM</p>	 <p>MON, 25 OCT 2021 9:00 AM</p> <p>International PHP</p>	 <p>MON, 08 NOV 2021 9:00 AM</p> <p>GOTO</p>	 <p>MON, 08 NOV 2021 8:00 AM</p> <p>W-IAx 2021</p>
--	---	--	--	--

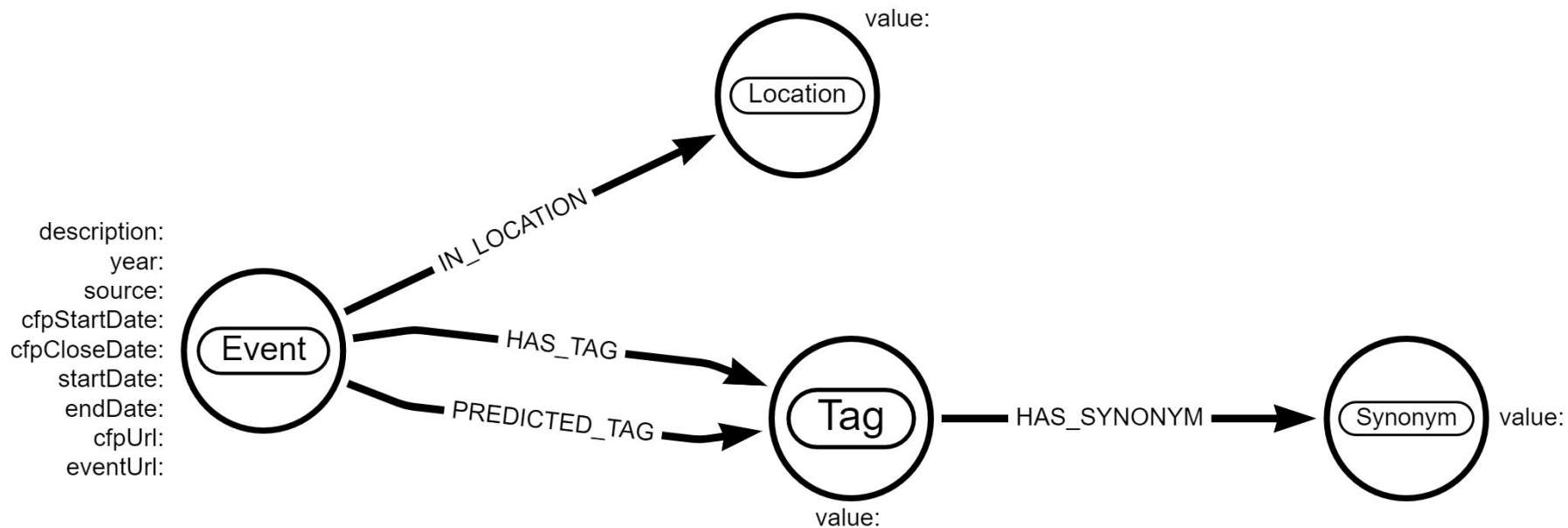
Building the CfP tool

1. Identify common places for Call for Papers information
 - a. Web search
 - b. Ask around
 - c. Past experience
2. Identify what data do we need and build a data model



The data model

Drawn using Arrows: <http://arrows.app>






Building the CfP tool

1. Identify common places for Call for Papers information
 - a. Web search
 - b. Ask around
 - c. Past experience
2. Identify what data do we need and build a data model
3. Figure out how to get at individual Call for Papers information
 - a. Is there a directory we can iterate over?
 - b. Is there a site map we can use?
 - c. Can we carefully craft a search engine query?
4. Examine the CfP source code to figure out how to extract the data
5. Extract the data and load into the database, based on the data model


[Home](#)
[Event Directory](#)
[Past Events](#)
[Pricing](#)
[Help](#)
[Login](#)


The image is a screenshot of a web application titled "Event Directory". At the top, there is a dark blue header with the title in white. Below the title, a navigation bar contains links for "All Events", "Open CFPs Only", and "Only Recurring Events/Meetups". On the right side of the header, there are two orange buttons: "Add Your Event!" and "Create Your CFP!". Below the header, the main content area has a light gray background. On the left, there is a sidebar with a checkbox for "Offers Travel Assistance" and a green button labeled "Save Search as an Alert". In the center, there is a search bar with the placeholder text "Search" and a magnifying glass icon. Below the search bar, a tagline reads "Tag search only: tags: first tag, second tag, etc...". The main content area displays a search result for "INTENT - The security research summit. - Virtual". The result card has a dark blue header with the event title and a share icon. Below the header, the text "INTENT" is followed by a URL "http://www.securityresearchsummit.com/". Further down, it lists "Event Dates: November 16, 2021" and "CFP closes at October 04, 2021 19:59 UTC" with a smaller note "October 04, 2021 20:59 BST (Local)". To the right of the event card, there are three green buttons: "Submit Now!", "I'm Attending!", and "Add to Calendar". Below the event card, there is a dark blue bar with the word "PRICING" and a dropdown arrow, followed by a blue button "CREATE EVENT", a link "LOGIN", and a button "REGISTER". Below this bar, there are two tabs: "ONLINE" and "PAST". At the bottom, there is a search bar with the placeholder text "Enter location" and a dropdown menu set to "All times", followed by a blue button "SEARCH". On the far right, there is a vertical sidebar with three green buttons: "Submit Now!", "I'm Attending!", and "Add to Calendar". At the bottom right corner, there is a green circular button with a white question mark and the word "Help".



MON, 11 OCT 2021 1:00 PM

Color Psychology for Behavioral & UX Design (Fall 2021)

📍 Online



WED, 13 OCT 2021 10:00 AM

Greece Gambling Conference 2021

📍 NJV Athens Plaza

```
157 <div class="justifize__box">
158   <div class="subheader__logo">
159     
160   </div>
161   <div class="subheader__group">
162     <h1 class="subheader__title">DjangoCon US 2021</h1>
163     <h1 class="subheader__subtitle">
164       Online
165       October 21, 2021, October 22, 2021, October 23, 2021
166     </h1>
167     <a target="_blank" href="https://2021.djangocon.us/">https://2021.djangocon.us/</a><br>
168     <span>Tags: <a href="/events?keywords=tags%3A+Python">Python</a>, <a href="/events?keywords=tags%3A+Django">Django</a>, <a href="/events?keywords=tags%3A+documentation">Documentation</div>
169   </div>
170
171 </div>
172 <div class="justifize__box pull-right">
173   <div class="subheader__subtitle">
174     <a href="https://www.facebook.com/dialog/feed?
175     app_id=929684357128596
176     &display=page&caption=PaperCall.io
177     &name=DjangoCon US 2021
178     &picture=https://www.papercall.io/assets/logo-papercall.svg
179     &link=https://www.papercall.io/djangocon-us-2021
180     &redirect_uri=https://www.papercall.io/djangocon-us-2021"><i class="fa fa-facebook" data-toggle="tooltip" data-placement="bottom" title="Share on Facebook"></i></a>
181     <a href="https://twitter.com/intent/tweet?text=Submit to the DjangoCon US 2021CFP! https://www.papercall.io/djangocon-us-2021"><i class="fa fa-twitter" data-toggle="tooltip" data-pla
182
183     <button class="fa fa-external-link copy-to-clipboard" data-clipboard-text="https://www.papercall.io/djangocon-us-2021" data-toggle="tooltip" data-placement="bottom" title="Copy to Cli
184   </h1>
185 </div>
186 </div>
187 </div>
188
189 </div>
190 </div>
191 </div>
192
193 <div class="container">
194   <div id="flash_notices">
195
196 <div class="row">
197   <div class="col-md-12">
```

Success! We're finished!



Not quite... We have some issues to deal with:

- Not all of the conference platforms that we looked at have tags
 - Use the tags we have
 - Scan across event descriptions and titles
- There are some data quality challenges
 - We'll look at some small fixes now
 - Explore options to sort in the next iteration

Are we done now?

Almost! We don't have many tags:

- Good source of technology tags - StackOverflow
- Use the StackExchange data explorer to pull top 200 tags
 - Also synonyms
- Some fuzzy matching after

Yet (more) issues:

- Not all tags make sense
- Use some basic stats to get rid of them
 - Proportion of talks with that tag
 - Frequency of tag appearing across titles/descriptions

Still no perfect, but good enough!

The tools in use

- Lots of googling
- Jupyter notebook
- Beautiful Soup & Google search
- StackExchange data explorer
- Neo4j stack
 - Sandbox - No-download trial database
 - Python driver - API to connect to the database
 - APOC library - a collection helper functions and procedures
 - Browser - developer aide for queries and visualisation
 - GDS library - graph algorithms for data science

**We're there! Let's find some
conferences to submit to!**



Let's go find some conferences!

- What are the most 'popular' tags?
- What CfPs are closing within the next month?
- What CfPs have the theme of data science/machine learning?
- What CfPs have tags of data science and python?
- What conferences can be grouped together by similarity?

So what's next?



What's next?

- General code and data clean-up
- More data!
- Find other sources for more tags and synonyms
- Use of NLP to extract more meaningful tags/exclude meaningless tags
- Create conference 'themes', e.g. DevOps, Data Science, etc.
- Load up historical conferences
- Explore options for importing historical talk titles

Repo: <https://github.com/lju-lazarevic/cfptool>

Want to learn more about graphs?

Free online training and certification:

- dev.neo4j.com/learn

How to, best practices, hands on and community stories:

- dev.neo4j.com/videos

Come say hello :)

- dev.neo4j.com/chat
- dev.neo4j.com/forum

