

Japanese Visual Media Graph: Project Update

Prof. Magnus Pfeffer, Stuttgart Media University
<pfeffer@hdm-stuttgart.de>

- Project starting point
- Lessons and results from project phase 1 (2019-2022)
- Project phase 2 (2023-2025)

Project starting point

- “Manga, Anime and Games” are a growing field of study
 - Interest from Japanese and Media Studies
 - Need for data-driven research options
- Lack of comprehensive resources
 - Rapidly growing number of media
 - Library catalogs and similar publisher-sourced data
 - Small encyclopedic efforts
- Vibrant enthusiast communities
 - Active since the late 1990s
 - Collect and curate information on Japanese Visual Media

- Create a research database on Japanese visual media, including, but not limited to, anime, manga, computer games and visual novels
- Aimed at researchers in Japan Studies who focus on modern media and its expressions, themes, topics, characters and reception
- Use a graph-based, highly interconnected database structure, similar to the Google knowledge graph, that is combined with a flexible search interface and analytic tools
- Use the data on Japanese visual media that is being created and curated by the many enthusiast communities on the web

Lessons and results from project phase 1 (2019-2022)

- Can be hard to contact
 - Emails not answered
 - Non-standard communication (IRC, Discord, Social Media)
- Need time to establish a connection
 - Are we even real? Or a troll on the internet?
 - Need to explain our academic motivations
- Very supportive once personal contact was made
 - In-person workshops with researchers were helpful

- The visual novel database (vnadb.net)
 - International community
 - Focussed on the visual novel game genre
- AnimeClick (animeclick.it)
 - Italian community
 - Broad interest in Japanese media and culture
- Anime Characters Database (animecharactersdatabase.com)
 - International community
 - Focussed on character appearance and personality
- Media-Arts Database
 - Japanese project, utilizing publisher data

- Heterogenous licensing by the communities
 - Low awareness of copyright in database context
 - Non-standard license agreements and site policies
 - Incompatible licenses
- Legitimate concerns of the communities
 - Wholesale copying of sites to divert traffic and profit by advertisements
 - Recognition of their work by proper attribution data sources
- License must cover all research use cases
 - Combining with other data, publication of analysis and research
- Solution: CC-BY-SA-NC 4.0, dual licensing where needed

- Extensive analysis of data accuracy
 - Based on title information (Japanese and English) of anime works
 - Comparison to “ground truth”: Media photos, OP title cards, publisher page
 - Sample size large enough to be in a 5% margin at 95% probability
- Results very encouraging
 - Very low error rate in enthusiast communities’ data
 - Most common type of error is typographical inconsistencies
 - Disambiguation information stored in “original title” field
- Wikidata is much worse in comparison
 - up to 10% misrepresentation errors, up to 16% missing data

- 1. Merge works based on title information (in Japanese kanji)
 - Use very strict string comparison first
 - Relax string comparison to address typographic inconsistencies next
 - Resolve ambiguous matches manually
- 2. Compare information on creators per work cluster
 - Match candidates can be found with relaxed name matching
- 3. Use creator clusters to find new candidates for matching works
- 4. Use n-gram comparison for remaining unmatched works
 - Manual resolution required

- Done
 - Anime and manga work clusters using strict and relaxed matching
 - Manual resolution of ambiguous matches of work titles of anime
 - Matching between multiple sources

- Ongoing
 - Creator clusters using Japanese names
 - Problem: creator data in media-arts database not normalized

- Future
 - Visual character clusters

- Idea
 - Showcase research using the knowledge graph
 - Limited scope to reach results more quickly
 - Used to direct the development in the project
- Examples
 - Census of characters in Japanese visual media
 - Testing a point of Hiroki Azuma's "Otaku: Japan's Database Animals"
 - Recurring patterns in character creation in visual novels

- Testing of different approaches
 - RDF triple store with separate web frontend and external search index
 - Integrated RDF/linked data hosting solution
 - Integrated graph database solution
- Performance issues
 - Viewing large result sets in web frontend
- Licensing issues
 - Proprietary licenses with unusual restrictions

- Integrated solutions not a good fit for our data and use cases
- Own solution
 - Apache Fuseki RDF triple store with SPARQL interface
 - Custom web frontend with optimized queries for large results
 - Elasticsearch index for search function
 - Open source



Tell us your
ideas!



<https://github.com/Japanese-Visual-Media-Graph>

Saiyuki Gaiden

Resource: <http://mediagraph.link/aclick/work/3222>

Graph: **aclick**

Property

Value

label

Saiyuki Gaiden

Category en

OVA series en Serie OAV it

English title en

Saiyuuki Gaiden

Episodes en

3

is From work en of

12

- Erlang Shen
- Jirōshin
- Kanzeon Bosatsu
- Kenren Taishō
- Konzen Dōji
- Li Tōten

- Nataku Taishi
- Seikai Ryūō Gōjun
- Seiten Taisei
- Shōu
- Son Goku
- Tenpō Gensui

Genre en

- Action en Azione it
- Adventure en Avventura it

- Fantastic en Fantastico it
- Supernatural en Soprannaturale it

ID en

3222

Kanji title en

最遊記外伝

Nationality en

Giappone

Dark Mode Search crosstab languages

1

3

en
it
x-jat
all

2

5 hide

4

3

- 1: CSS selection
- 2: Source indicator
- 3: Language filter
- 4: Expandable sets
- 5: Source filter

Dark Mode

JVMG-Search

dragonball Search

type

- Game package (12)
- Manga book (7)
- Manga book series (4)
- Game variation (1)

graph

- madb (24)

Number of results: 24

uri: <http://mediagraph.link/madb/id/M721299>
label:
• DRAGONBALL FighterZ

type:
• Game package

graph:
• madb

uri: <http://mediagraph.link/madb/id/M721298>
label:
• DRAGONBALL FighterZ

type:
• Game package

graph:
• madb

Search all indexes

Facets for data source
and entity type
(OR-combined)

Result list
with match highlight

Project phase 2 (2023-2025)

- Federico Pianzola (University of Groningen)
 - Fan works
- Martin Hennig (University of Tübingen)
 - Media genres
- Bryan Hartzheim (Waseda University Tokyo), Stevie Suan (Hosei University Tokyo)
 - Creator careers

A locally running software stack

- Server mediagraph.link has limited resources
- Local copy of database useful for larger analysis tasks
- Goal
 - Local copy using container or other technology
 - Easy to install and run on any operating system (Win, Mac, Linux)
 - Integration into local research environment (Jupyter Labs)

- Documentation
 - Complete documentation of all ontologies
 - Video guides and online tutorials for common tasks
 - Example complex analysis tasks using Python and Jupyter Labs
- Workshops
 - Single-day workshops at conferences
 - Multi-day interactive workshops at locations Stuttgart and Kyoto
 - Open to all interested researchers

Thank you for your attention.