

# Web Scraping and OSINT Tools for Crime Investigation: Detecting Suspicious Seller Networks

Presented By:  
T.U Adarsh-2023BCY0039.



# Objectives

The primary objective of this project is to demonstrate how web scraping techniques can be applied to support crime investigations by identifying patterns of potentially stolen goods on online marketplaces. A privately hosted test website containing synthetic product listings will be created to simulate a real marketplace environment. Using Python tools such as Requests, BeautifulSoup, or Selenium, data including product titles, descriptions, prices, seller details, contact information, image URLs, and timestamps will be collected and analyzed. The goal is to detect repeated indicators of suspicious activity — such as reused images, identical contact details, and repeated phrases — and visualize the relationships between sellers and listings through a network map. This will highlight clusters of interconnected or potentially illicit behavior. The project also aims to document the methodology, findings, and ethical considerations, confirming that all data used is synthetic and hosted privately for academic and demonstration purposes.

# Used Tools and modules

- Web Scraping (Python + BeautifulSoup)
  - Purpose: Extract product listings automatically from the synthetic marketplace website.
- Web Scraping (OctoParse)
  - Octoparse is a no-coding solution for web scraping to turn pages into structured data.
- Data Flagging & Cleaning (Pandas & OpenRefine)
  - Purpose: Clean, standardize, and flag suspicious attributes like reused contacts, images, and low prices.
- Network Visualization (Gephi & Plotly)
  - Purpose: Visualize relationships between sellers, contacts, and reused images to identify hidden clusters.
- Unsupervised Machine Learning (Algorithm-Isolation Forest)

Purpose: Predict and classify suspicious listings based on risk features for automated detection.
- OSINT Tool – Maltego
  - Purpose: Enrich investigation by mapping connections from emails, domains, and online profiles.
- OSINT Tool – Numverify
  - Purpose: Validate phone numbers to check authenticity and detect disposable or fake numbers.

# Web Scraping:

- Built using Python, leveraging Requests and BeautifulSoup libraries for efficient data extraction.
- Scraps synthetic product listings from a locally hosted website (ethical, no real data).
- Extracts key fields: title, description, price, seller handle, contact info, image URL, timestamp.
- Identifies repeated indicators such as reused images, duplicate contacts, and identical descriptions
- Cleans and stores the extracted data into a structured CSV dataset for analysis.
- Designed for speed, simplicity, and reusability, ensuring consistent results across runs.

A simple flow diagram:

```
Website → HTML Fetch (Requests) → Parse Data (BeautifulSoup) → Clean & Store  
(CSV)
```

## 1. Applications of Web Scraping

- Price monitoring and comparison (e-commerce)
- Sentiment analysis from social media
- Job listing aggregation
- Data collection for machine learning
- Academic or market research data gathering

## 3. Future Improvements

- Automating with a scheduler (cron job)
- Integrating NLP for analyzing scraped text
- Using cloud-based scraping (AWS Lambda, Google Cloud Functions)

## 2. Challenges & Limitations

- Dynamic websites (JavaScript-rendered content)
- Captchas and anti-bot mechanisms
- Rate limiting and IP blocking
- Legal and ethical boundaries (robots.txt, GDPR)

## 4. Advanced Tools

- Selenium / Playwright – handles JavaScript-heavy sites
- API scraping – more reliable than HTML parsing if available
- octoparse- no code web scraping tool

# octoparse tool:

You can open any website inside its built-in browser, click on the elements you want (like product names, prices, emails, etc.), and it will automatically create a scraping workflow.

You can then run the scraper locally on your computer or in the cloud to extract the data into formats like CSV, Excel, JSON, or send it directly to a database or API.

The image shows the Octoparse software interface. On the left, a 'Task List' sidebar includes icons for Home, Task List, New Task, Up..., + New, Task List, Templates, Tools, Pricing, RPA, Referral, Inbox, Help, and Settings. A main window titled 'DarkMarket - Stolen Go...' displays a 'Product Listings' page with three items: 'Rolex Submariner Watch', 'Diamond Necklace', and 'Omega Speedmaster Watch'. A central modal window titled '+ Tips' says 'List data detected. Preview your data at the bottom and choose how you'd like to extract the data.' It has a checkbox 'Extract data list' and buttons 'Create workflow' and 'Cancel'. Below the modal is a 'Data Preview' table with 7 rows of sample data. The table has columns: No., Pre, Title, Image, Description, Price, sellerhandle, Info, Time, +, and Actions. The first row shows 'Rolex Submariner W...' with 'Pre' checked. The second row shows 'Diamond Necklace'. The third row shows 'Omega Speedmaster...'. The fourth row shows 'Gold Engagement Ri...'. The fifth row shows 'iPhone 14 Pro'. The sixth row shows 'MacBook Pro 16"'. The seventh row shows 'Dell XPS 15'. The right side of the interface shows a summary of a completed task: '165 Data Extracted' (circled in green), 'Completed', 'Task completed', 'Duplicates: 0 lines', 'Time Spent: 9s', 'Avg. Speed: 1016 lines/min', 'Rerun' button, and 'Export' button. Below this is a 'Data List' table with 11 rows of extracted data, matching the preview in the center. The table columns are: #, Title, Image, Description, Price, sellerhandle, Info, and Time. The first row shows '1 Emerald Halo Necklace'. The second row shows '2 Patek Philippe – Collect...'. The third row shows '3 iPhone 14 – Dual SIM'. The fourth row shows '4 Samsung – Promo Unit'. The fifth row shows '5 Chromebook – Lightwei...'. The sixth row shows '6 Rolex – Fine Timing'. The seventh row shows '7 Diamond Pendant – Cla...'. The eighth row shows '8 Patek Philippe – Invest...'. The ninth row shows '9 iPhone 14 – Family Bun...'. The tenth row shows '10 Samsung – Excellent Co...'. The eleventh row shows '11 Used Gaming Laptop – ...'. The table includes pagination controls at the bottom.

#	Title	Image	Description	Price	sellerhandle	Info	Time
1	Emerald Halo Necklace	http://127.0.0.1:5500/fi...	Beautiful emerald cente...	\$3900	JewelThief	jewels@fake.com	2024-01-29
2	Patek Philippe – Collect...	http://127.0.0.1:5500/fi...	Vault release for collect...	\$45000	WatchDealer99	watchdealer@example....	2024-01-30
3	iPhone 14 – Dual SIM	http://127.0.0.1:5500/fi...	Works with dual SIMs, f...	\$800	TechAlias	techguy@alias.net	2024-01-31
4	Samsung – Promo Unit	http://127.0.0.1:5500/fi...	Promotional unit in exc...	\$600	MobileMafia	contact@mobilemafia.net	2024-02-01
5	Chromebook – Lightwei...	http://127.0.0.1:5500/fi...	Perfect for students and...	\$260	BagBroker	broker@bags.example	2024-02-02
6	Rolex – Fine Timing	http://127.0.0.1:5500/fi...	A fine selection of Rolex...	\$20000	WatchDealer99	watchdealer@example....	2024-02-03
7	Diamond Pendant – Cla...	http://127.0.0.1:5500/fi...	Elegant piece, perfect c...	\$1600	JewelThief	jewels@fake.com	2024-02-04
8	Patek Philippe – Invest...	http://127.0.0.1:5500/fi...	Investment-grade time...	\$42000	WatchDealer99	watchdealer@example....	2024-02-05
9	iPhone 14 – Family Bun...	http://127.0.0.1:5500/fi...	Bundle offer for multipl...	\$3000	TechAlias	techguy@alias.net	2024-02-06
10	Samsung – Excellent Co...	http://127.0.0.1:5500/fi...	Seller-rated high trust, c...	\$630	MobileMafia	contact@mobilemafia.net	2024-02-07
11	Used Gaming Laptop – ...	http://127.0.0.1:5500/fi...	Strong CPU and GPU co...	\$1600	GadgetGhost	ghost@gadget.net	2024-02-08

# Red-Flag Indicators

To detect patterns associated with stolen or fraudulent goods, eight red-flag indicators were engineered from the scraped marketplace data.

These features captured both behavioral and content-based irregularities:

- Reused Image – Identifies listings that use the same product image across multiple sellers, indicating duplicate or fake listings.
- Shared Contact – Detects repeated contact details (emails or phone numbers) across different sellers, hinting at coordinated or false identities.
- Suspicious Phrase – Flags risky text such as “urgent sale,” “no bill,” “pickup only” that may imply hidden ownership or unverifiable goods.
- Price Outlier – Captures abnormally low-priced items that fall below 50% of the typical product median, signaling possible stolen goods.
- Seller Burst (1-hour / 6-hour) – Highlights sellers posting several listings within short time spans, a sign of automation or bulk dumping.
- High-Velocity Seller – Identifies sellers posting at unusually high daily frequencies, suggesting organized reselling activity.
- Near-Duplicate Description – Detects copied or lightly modified listing descriptions that often appear in fraudulent or spammy posts.

# OpenRefine tool

- OpenRefine is a free, open-source tool for cleaning and organizing messy data.
- It works like a spreadsheet but is more powerful for data transformation.
- Helps detect and fix inconsistencies, duplicates, and formatting errors.
- Supports data import and export from various formats like CSV, Excel, JSON, and databases.
- Commonly used for data preparation before analysis or machine learning tasks.
- It allows advanced data filtering and clustering to quickly find hidden patterns or errors.
- Can connect to web APIs for fetching or reconciling data with external sources.
- Uses powerful **scripting languages (GREL, Python, or Clojure)** for custom data transformations.

## Limitations:-

- OpenRefine cannot handle very large datasets efficiently due to memory limitations.
- It lacks real-time collaboration or multi-user editing features.
- Some advanced operations require knowledge of GREL or scripting, which may be difficult for beginners.

# OpenRefine tool

OpenRefine labelled dataset final csv [Permalink](#) [Open...](#) [Export](#) [Help](#)

Facet / Filter Undo / Redo 0 / 0 < 104 rows Extensions Wikibase

Show as: rows records Show: 5 10 25 50 100 500 1000 rows « first < previous 1 - 10 next » last »

reused_image	shared_contact	suspicious_phrase	price_outlier	mins_since_prev_by_seller	seller_burst_60m	seller_burst_6h	posts_per_day	high_velocitySeller	desc_clean	dup_description_near	risk_score_weighted
True	False	True	False		False	False	1.0	False	must sell fast, price slightly negotiable.	True	4.5
True	False	False	False		False	False	1.0	False	selling as i upgraded to a newer model.	True	2.5
True	False	True	False		False	False	1.0	False	urgent sale, need cash fast.	True	4.5
True	False	True	False		False	False	1.0	False	pickup only, no returns.	True	4.5
True	False	False	False		False	False	1.0	False	well maintained and tested before listing.	True	2.5
True	False	True	False		False	False	1.0	False	pickup only, no returns.	True	4.5
True	False	False	False		False	False	1.0	False	selling as i upgraded to a newer model.	True	2.5
True	False	True	False		False	False	1.0	False	no bill or receipt available.	True	4.5
True	False	True	False		False	False	1.0	False	no bill or receipt available.	True	4.5
True	False	False	False		False	False	1.0	False	slight scratches, otherwise in great condition.	True	2.5

104 rows					
Show as:		Show: 5 10 25 50 100 500 1000 rows			
	suspicious_phrase	price_outlier	mins_since_prev_by_seller	seller_burst_60m	seller_burst_6h
True	Facet ►			False	False 1.0
False	Text filter			False	False 1.0
True	Edit cells ►			False	False 1.0
True	Edit column ►		Split into several columns...		False 1.0
True	Transpose ►		Join columns...		False 1.0
False	Sort...		Add column based on this column...		False 1.0
True	View ►		Add column by fetching URLs...		False 1.0
False	Reconcile ►		Add columns from reconciled values...		False 1.0
True	False		Rename column...		False 1.0
True	False		Remove this column		False 1.0
False	False		Move column to beginning		False 1.0
			Move column to end		
			Move column left		
			Move column right		

# OpenRefine tool

Add column based on column price\_num

New column name

On error  set to blank  store error  copy value from original column

Expression  General Refine Expression Language (GREL)

```
value
```

No syntax error.

Preview History Starred Help

row	value	value
1.	77016.0	77016.0
2.	14235.0	14235.0
3.	122892.0	122892.0
4.	94597.0	94597.0
5.	141328.0	141328.0
6.	120792.0	120792.0

OK Cancel

Add column based on column price\_num

New column name  price\_flag

On error  set to blank  store error  copy value from original column

Expression  General Refine Expression Language (GREL)

```
if(value > 15000.0, true, false)
```

General Refine Expression Language (GREL) syntax error.

Python / Jython

Clojure

Preview History Starred

row	value	if(value > 15000.0, true, fals ...
1.	77016.0	false
2.	14235.0	false
3.	122892.0	false
4.	94597.0	false
5.	141328.0	false
6.	120792.0	false

OK Cancel

# Weighted Risk Score Calculation

To quantify the overall suspiciousness of each listing, a weighted risk-scoring system was designed.

This converts the eight red-flag indicators into a single interpretable metric that measures risk severity.

- Purpose: Combine multiple behavioral and textual warning signs into one numeric score to rank listings by their likelihood of being fraudulent.
- Weighted Approach:
  - The first four indicators (Reused Image, Shared Contact, Suspicious Phrase, Price Outlier) were given double weight ( $2\times$ ) because they directly suggest fraudulent intent.
  - The behavioral indicators (Seller Burst 1h/6h, High-Velocity Seller) were given medium weight ( $1\times$ ) to represent suspicious posting activity.
  - The Near-Duplicate Description indicator was given light weight ( $0.5\times$ ) since repetition alone may not imply fraud.
- Calculation Formula:
- Each listing's Risk Score =  $\Sigma$  (weight  $\times$  indicator value).
- Example →
- Risk Score =  $2*(reused\_image + shared\_contact + suspicious\_phrase + price\_outlier) + 1*(seller\_burst\_60m + seller\_burst\_6h + high\_velocity\_seller) + 0.5*(dup\_description\_near)$
- Threshold Selection:
- Listings with a weighted score  $\geq 3$  were labeled Suspicious (1), while others were considered Normal (0).
- This threshold balanced detection sensitivity and false positives.
- Risk Level Classification:
  - Low Risk: Score  $< 2$
  - Medium Risk:  $2 \leq \text{Score} < 3$
  - High Risk: Score  $\geq 3$

# GEPHI TOOL FOR VISUALIZATION

## Building the Network Map.gxep

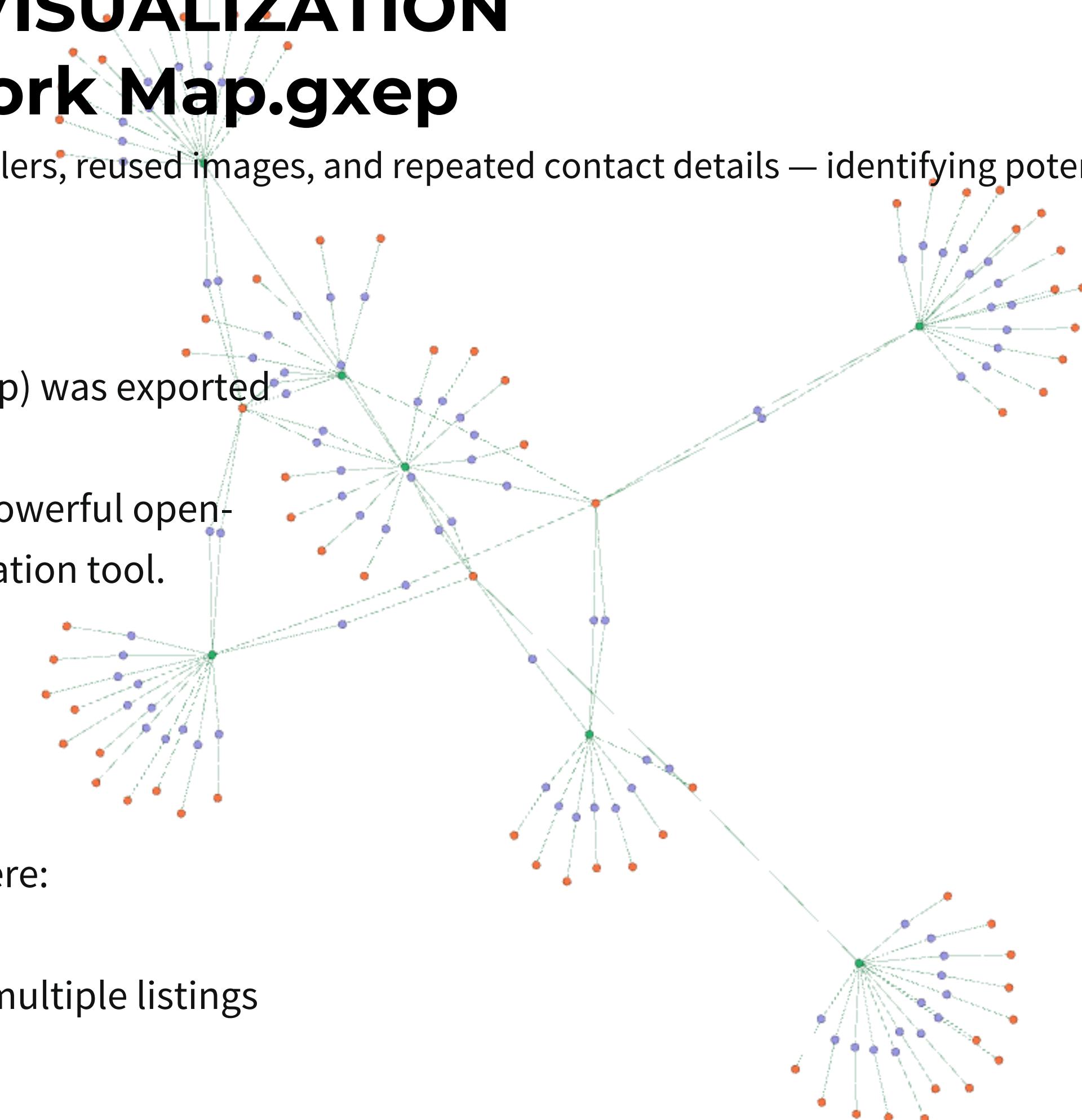
To uncover hidden connections between sellers, reused images, and repeated contact details — identifying potential illicit activity clusters on the marketplace.

### Methodology

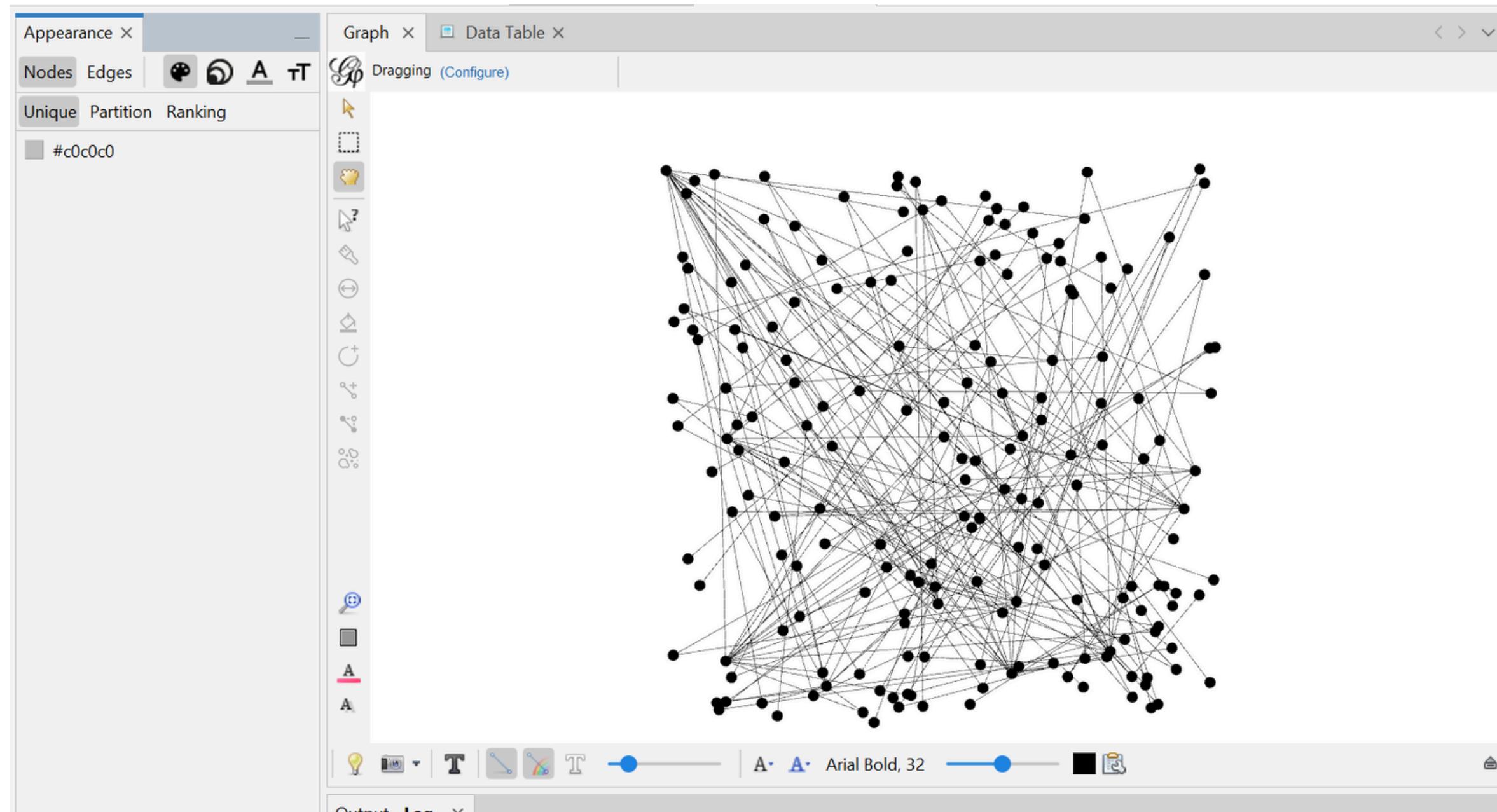
- The scraped data (from BeautifulSoup) was exported as a .gexf file.
- The file was imported into Gephi, a powerful open-source network analysis and visualization tool.
- Nodes represented:

- Sellers
- Contact Numbers
- Image URLs

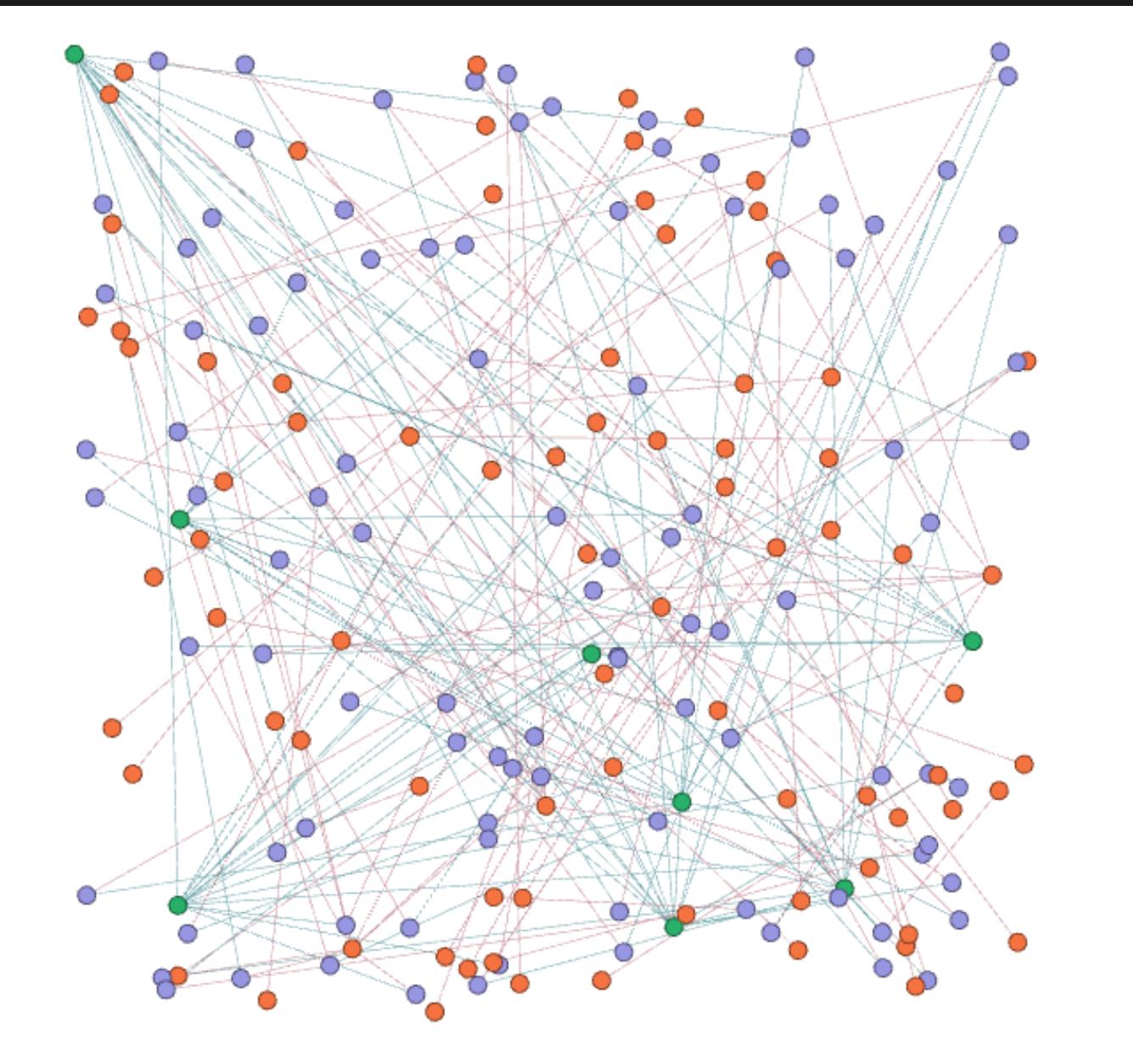
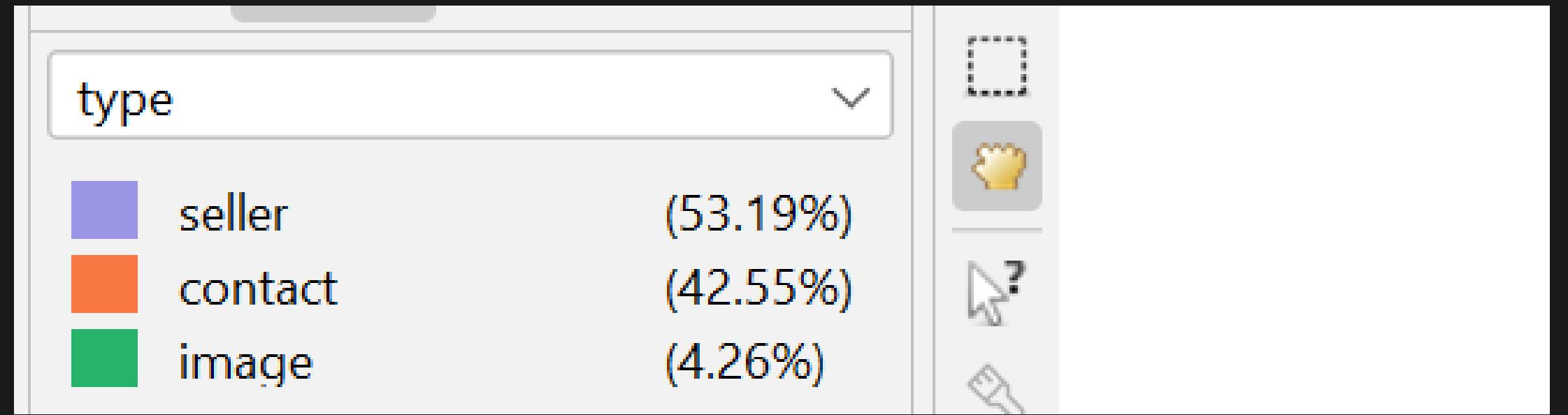
- Edges (connections) were drawn where:
  - The same image was reused
  - The same contact was linked to multiple listings



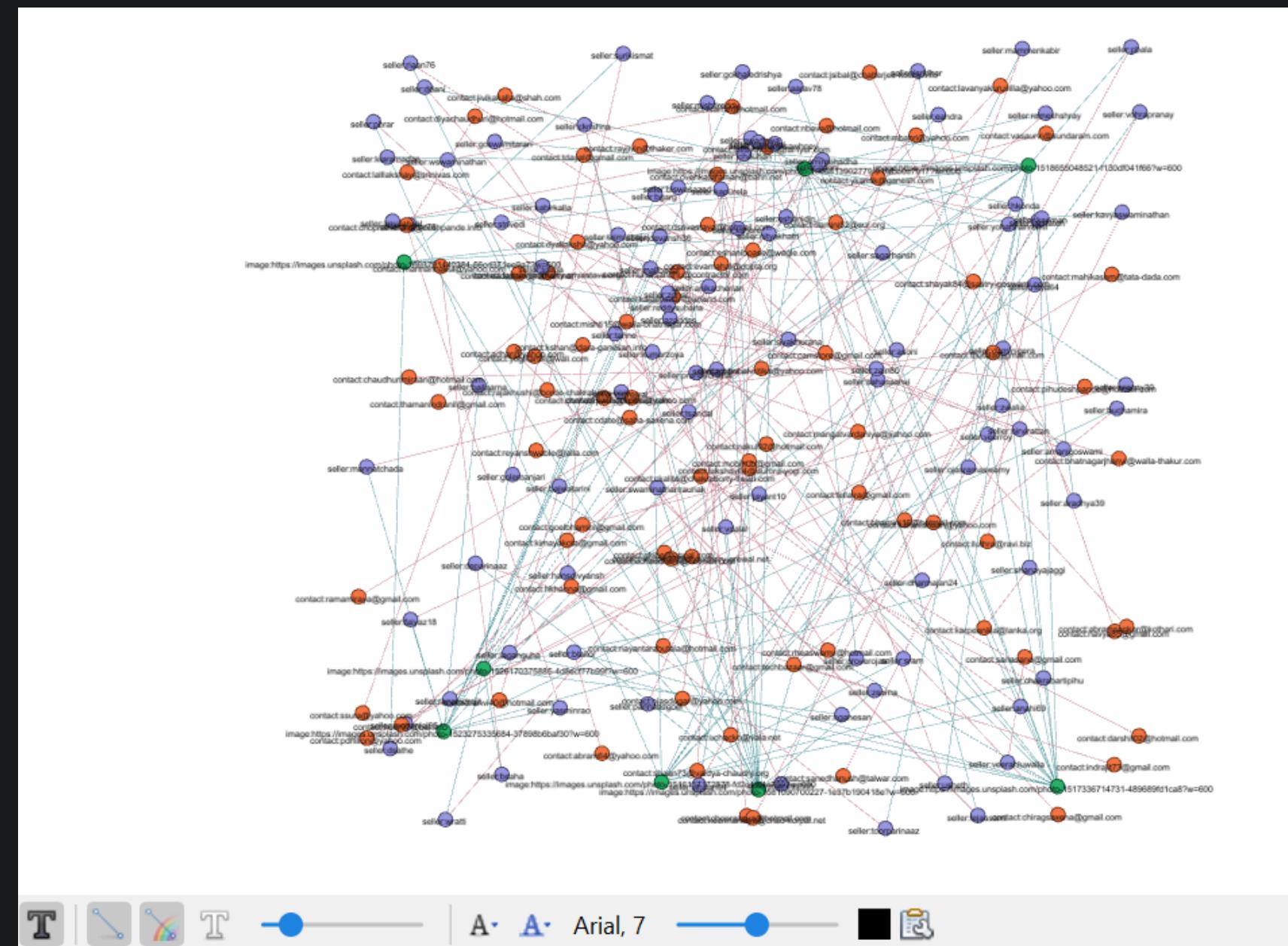
# Network Graph without Filters or Legend



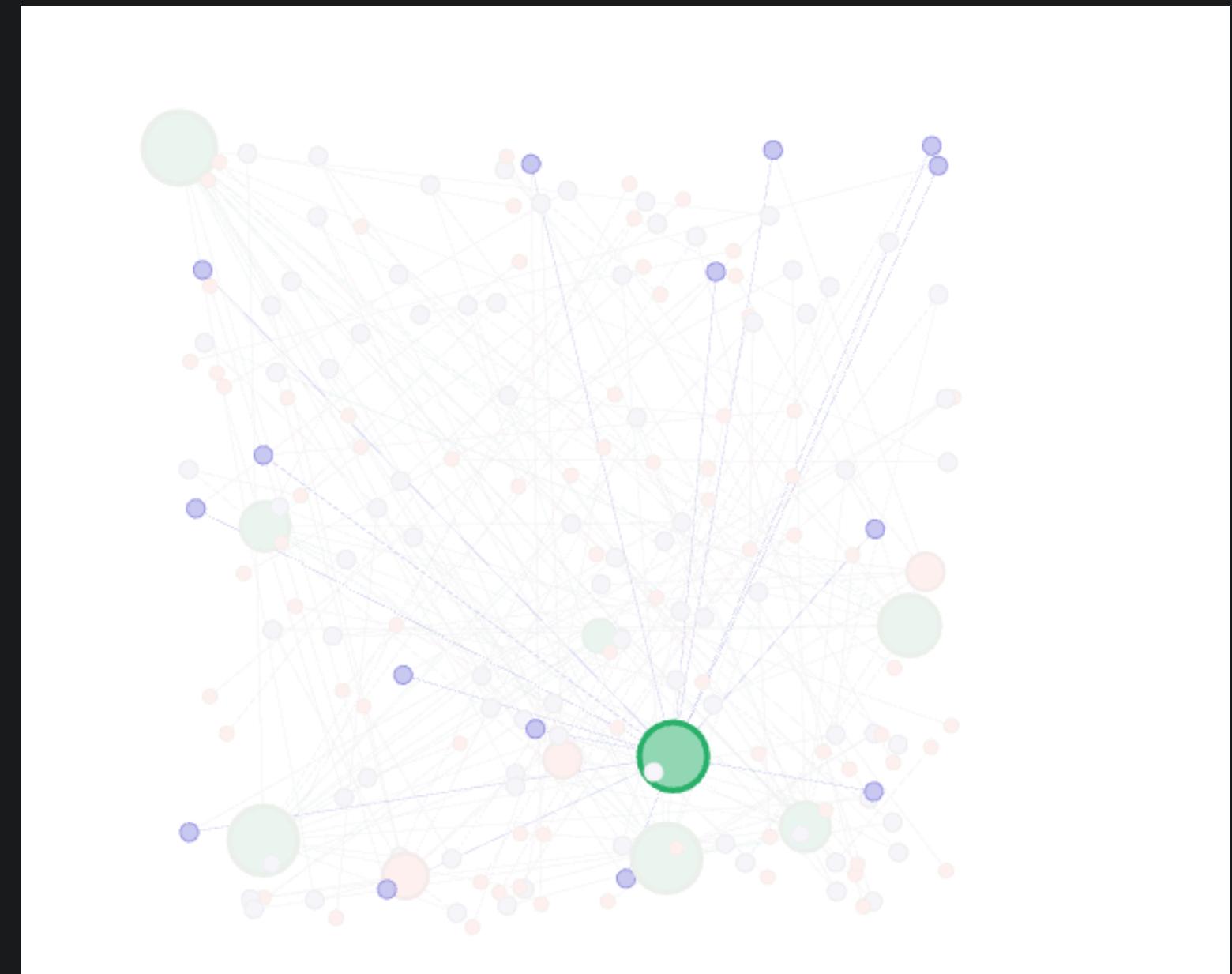
# Visualizing Relationships by Type: Sellers, Contacts, and Images



# ADDING LABELS TO NODES.

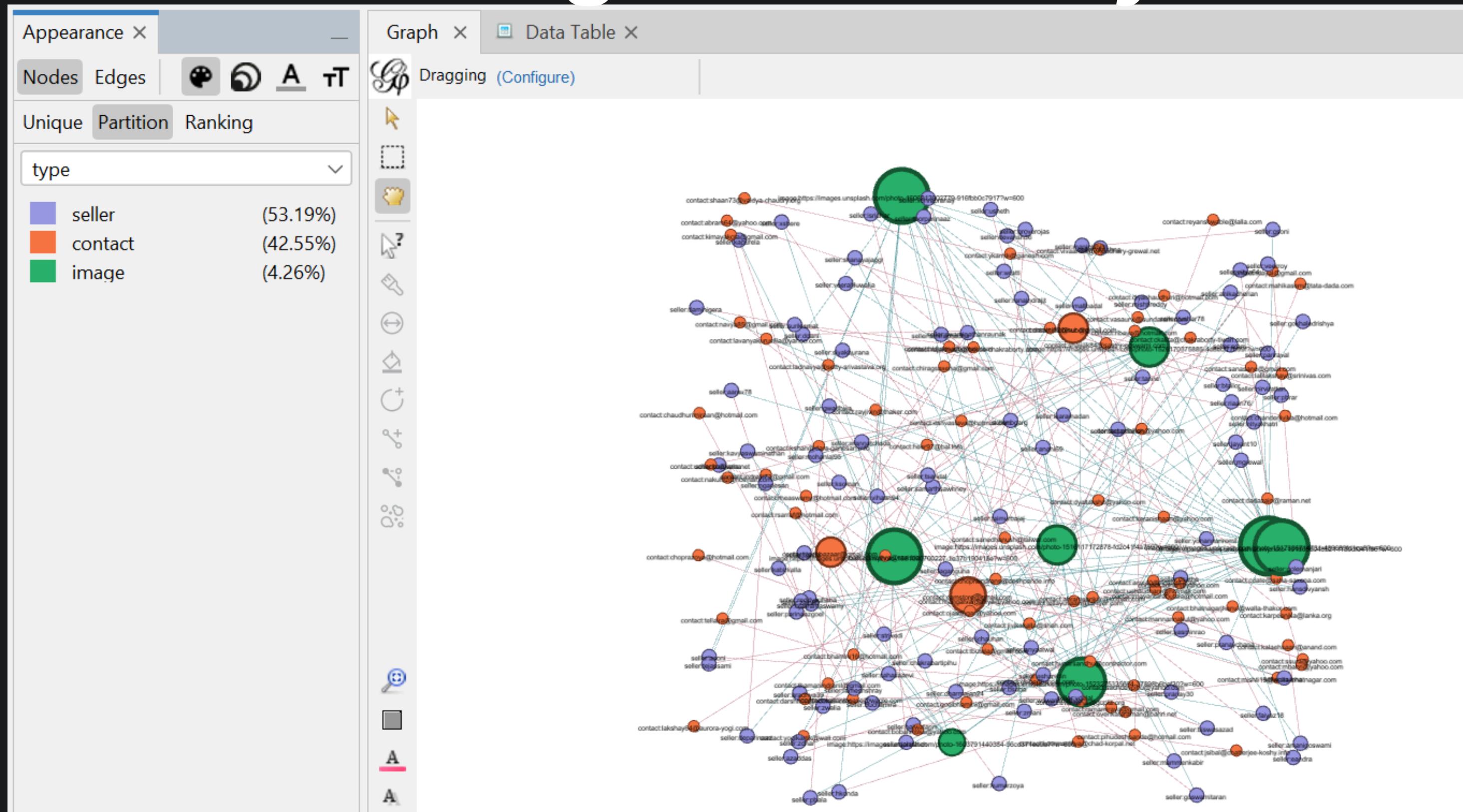


**Hover on a node to  
get its connections.**

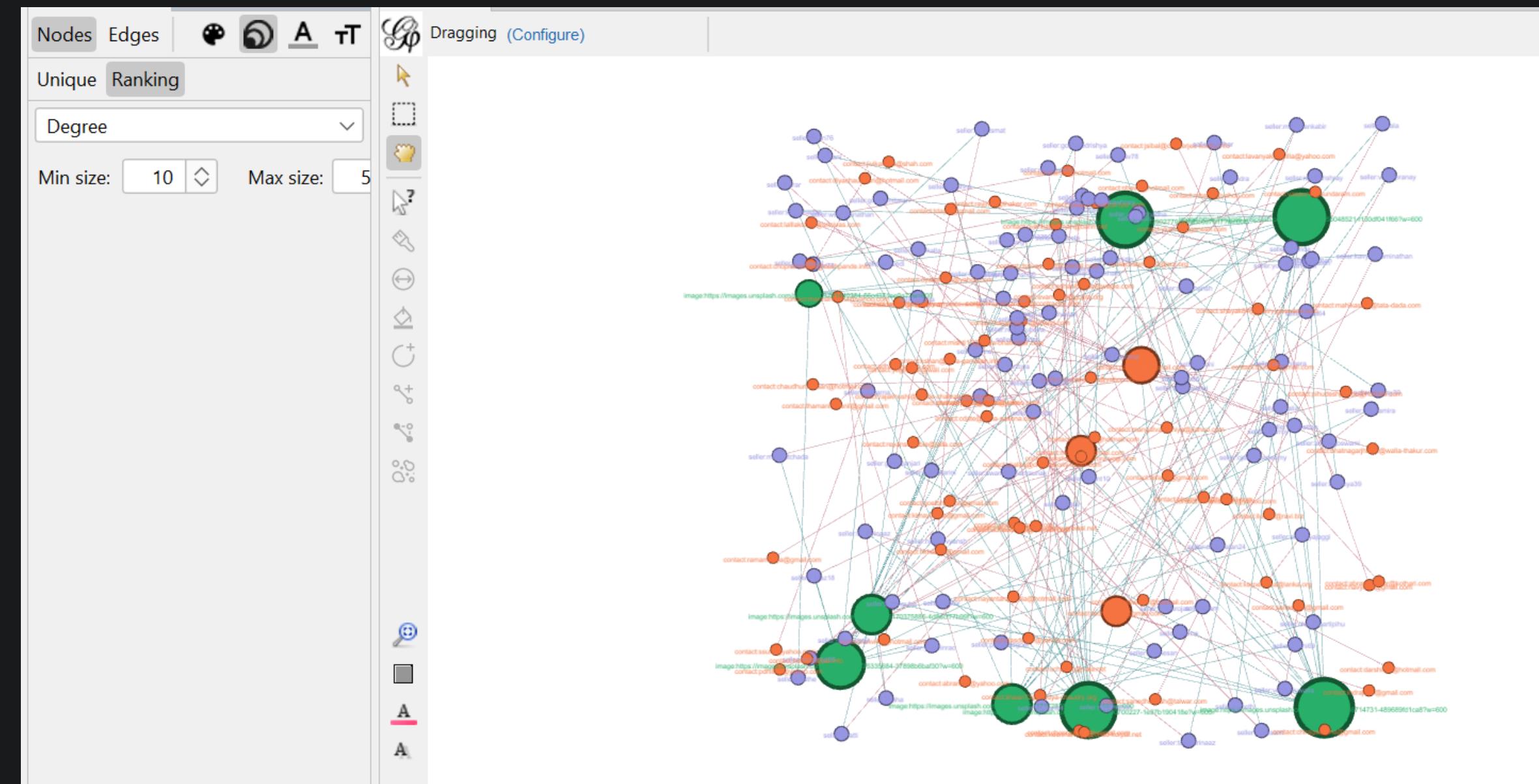


# seller\_network.gexf

(SIZE=degree centrality)

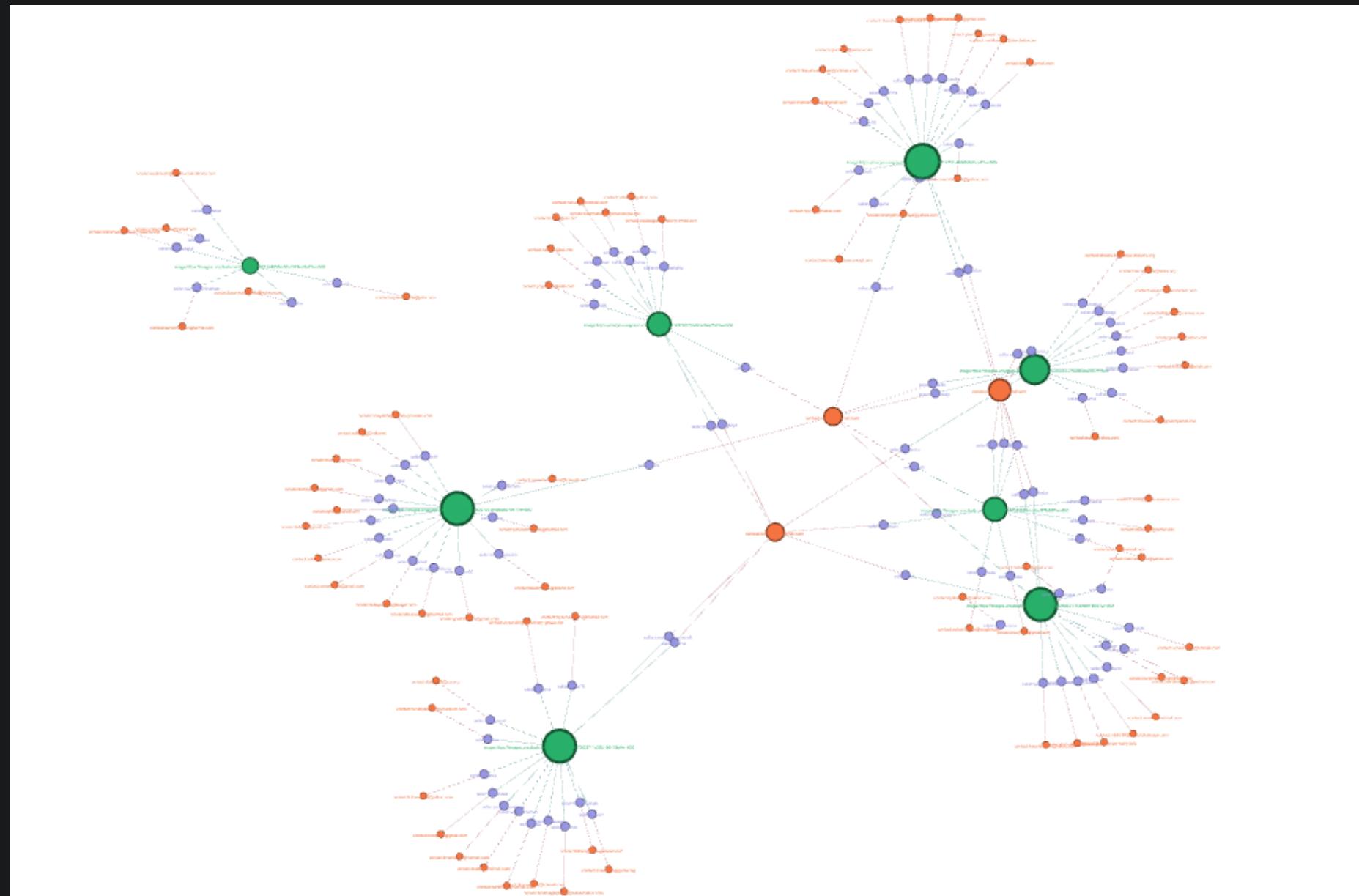


# LABEL COLOR IS MADE SAME AS NODE COLOR (BY CLICKING ON “A” IN TOP LEFT AND APPLY)

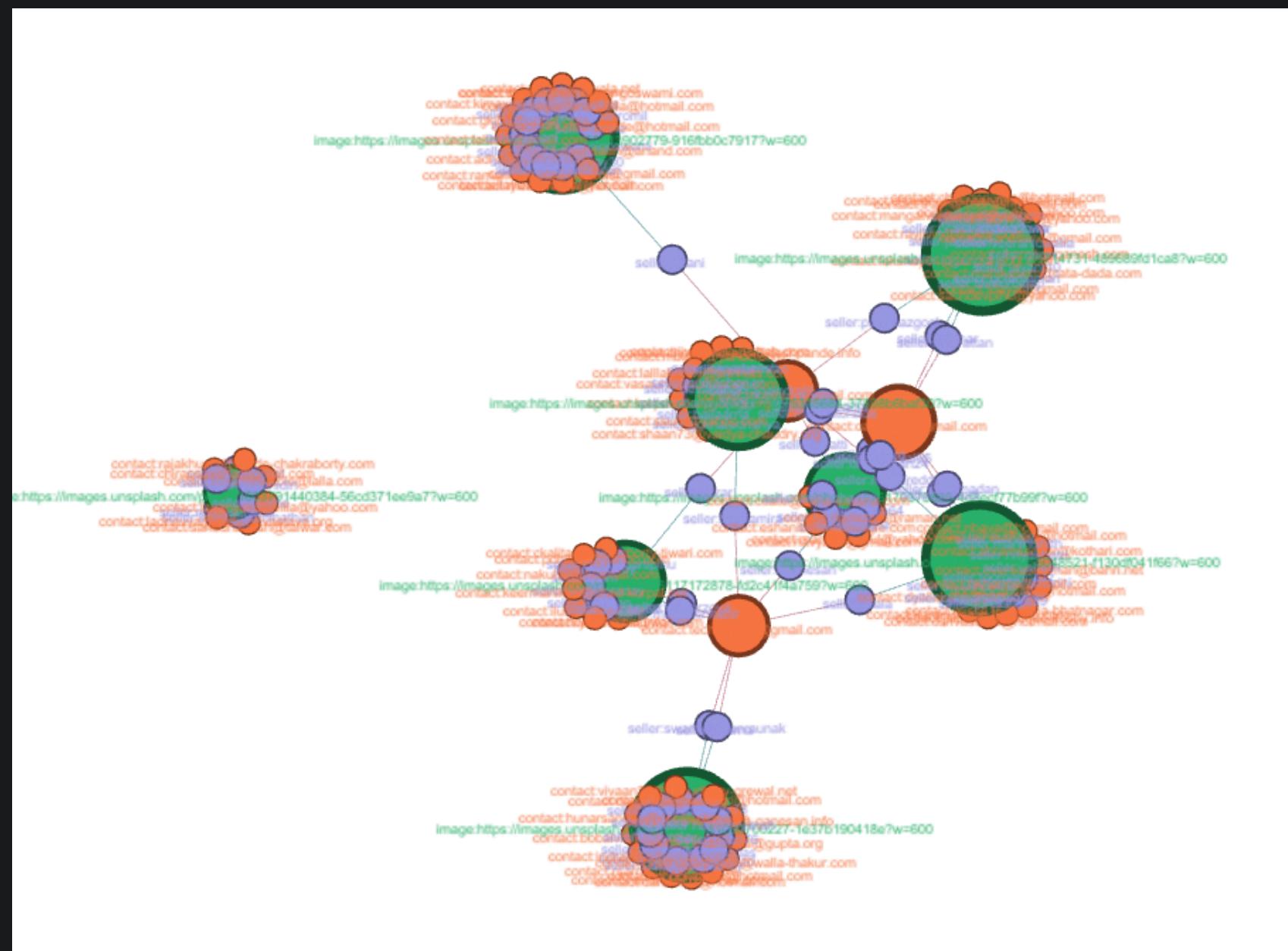


# EXPLORING DIFFERENT INBUILD LAYOUTS IN THE TOOL.

## YIFAN HU LAYOUT:



# FORCE ATLAS LAYOUT:



# Hybrid Clustering Approach for Suspicious Seller Detection

The objective of this module is to identify patterns of coordinated or fraudulent behaviour among online marketplace sellers. Instead of analysing each listing independently, the approach looks for relationships between listings — either through textual similarity or shared digital identifiers — to uncover hidden connections that may indicate duplicate or fake seller accounts

## ● Text-Based Clustering

focuses on the linguistic and semantic similarity of product descriptions.

Listings that describe products in nearly identical language are grouped together, even when the wording is slightly altered.

This helps expose reused or templated descriptions across different accounts.

## ● Graph-Based Clustering

models the relationships between sellers as a social-network graph.

Connections are created whenever two sellers share the same contact information or reuse the same product image, signalling possible coordination or duplicate identities.

**By combining textual and relational evidence, the system detects potential multi-account misuse, fake listings, and unauthorised reselling. It provides a structured, data-driven way to visualise and verify suspicious clusters that would be difficult to notice manually.**

# Text-Based Clustering (Semantic Similarity)

Text-based clustering aims to detect listings that use similar or identical product descriptions across different sellers.

Even when words or sentence structures vary, sellers may reuse the same underlying text patterns — often a strong indicator of coordinated copied behaviour.

## Methodology

- Product descriptions are first transformed into semantic representations using Sentence-BERT, a deep learning model trained to understand sentence meaning rather than just keywords.
- The similarity between every pair of listings is then measured using cosine similarity, which quantifies how closely their meanings align.
- Listings with similarity scores above a predefined threshold (e.g., 0.82) are grouped together into clusters, representing semantically identical or near-identical descriptions.
- Each cluster represents a potential case of cross-seller text reuse or copy-paste listings.

## Interpretation

- When multiple sellers use nearly identical descriptions for similar products, it may indicate:
  - Template-based fraudulent posting, where the same seller operates multiple accounts.
  - Content plagiarism, where legitimate sellers copy text from others to appear authentic.
  - Automated listing generation, where bots replicate product descriptions.

## Outcome:

- **The approach identified multiple semantic text clusters containing listings that are contextually similar despite minor rewording.**
- **These clusters highlight reused content patterns that can serve as an early warning for fraudulent or duplicate listings.**

# Text-Based Clustering-Output

Found 17 suspicious text clusters (cross-seller reuse)				
Cluster ID	Seller	Title	Description	
0	0	aarav78	HP Pavilion 15	Must sell fast, price slightly negotiable.
1	0	amirachadha	DJI Pocket 2	Urgent sale, need cash fast.
2	0	eandra	Sony Alpha a6400	Urgent sale, need cash fast.
3	1	amanigoswami	Xbox Series S	Selling as I upgraded to a newer model.
4	1	aradhya39	Royal Enfield Classic 350	Selling as I upgraded to a newer model.
5	1	bajwatarini	Yamaha MT-15	Selling as I upgraded to a newer model.
6	2	anahi69	Boat Rockerz 450	Pickup only, no returns.
7	2	anyaatwal	Steam Deck	Pickup only, no returns.
8	2	btailor	Dell XPS 13	Pickup only, no returns.
9	3	anikacherian	Fossil Chronograph	Well maintained and tested before listing.
10	3	ckrishna	Black+Decker Saw	Well maintained and tested before listing.
11	3	golemanjari	HP Pavilion 15	Well maintained and tested before listing.
12	6	bgarg	DJI Pocket 2	Used for a few months, works perfectly.
13	6	kabirkalla	Canon EOS 250D	Used for a few months, works perfectly.
14	6	maninishith	Honda Shine	Used for a few months, works perfectly.

## Observations

The text-based clustering grouped listings with semantically similar descriptions, revealing sellers who use almost identical product texts.

Each cluster represents potential reused or templated listings — for example, multiple sellers using the same sales pitch like “**Selling as I upgraded to a newer model**” or “**Pickup only, no returns**”.

This helps flag coordinated or duplicate seller behavior across the marketplace.

# Graph-Based Clustering (Seller Relationship Network)

Graph-based clustering extends the analysis beyond textual data by modelling relationships between sellers based on shared attributes. This method visualises how sellers are connected through shared digital identifiers, enabling the detection of groups that operate as coordinated entities.

## Methodology

- Each seller is represented as a node in a network graph.
- Edges (connections) are created whenever two sellers share:
  - The same contact information (e.g., phone number or email), or
  - The same product image reused across listings.
- These relationships form a web of interconnected sellers.
- Using a modularity-based community detection algorithm, the network is divided into communities (clusters) where sellers are densely interconnected.
- Each community thus represents a potential group of accounts linked through shared identifiers.

## Interpretation

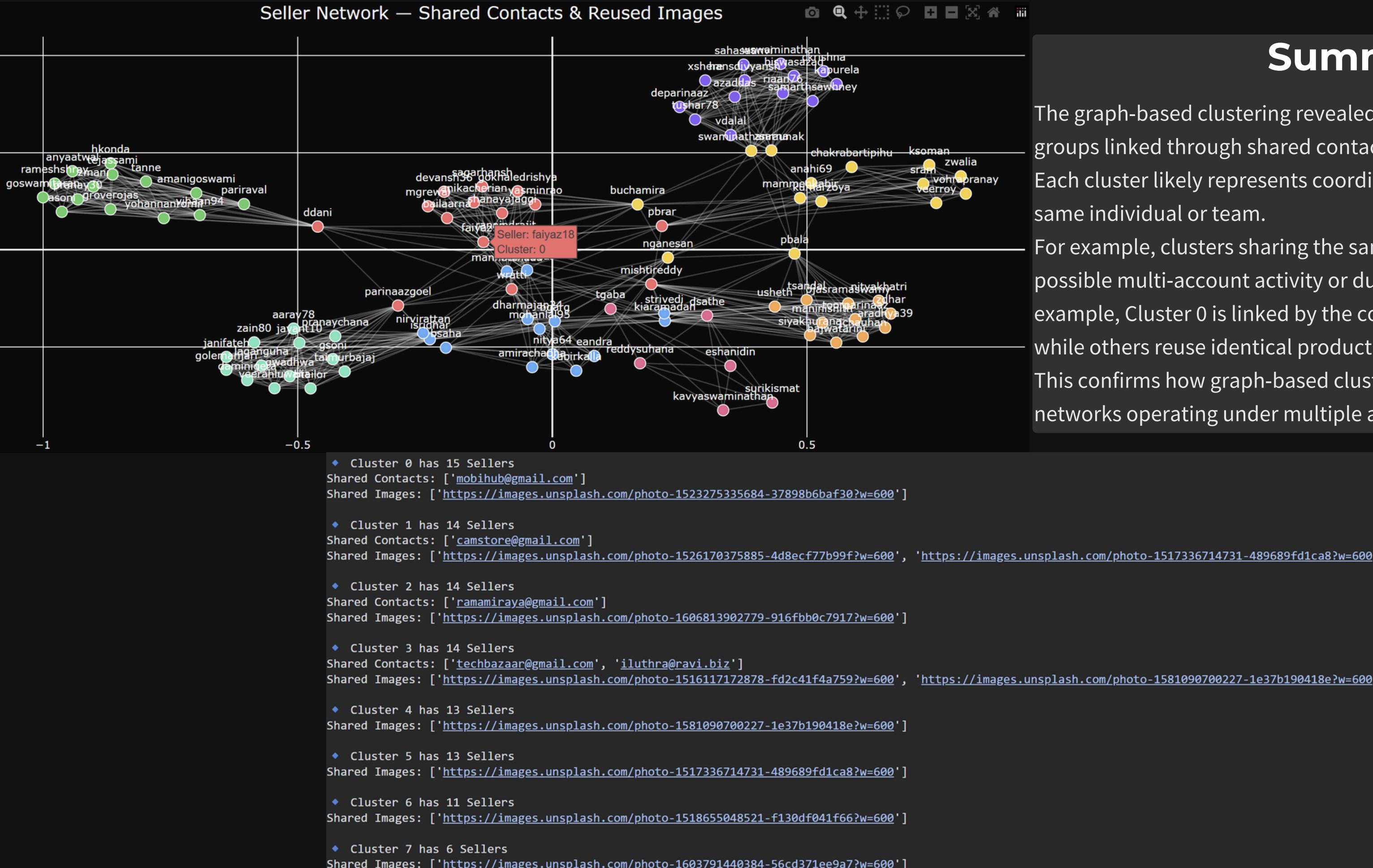
- Sellers within the same cluster are likely managed by the same individual or team, even if their public account names differ.
- Shared contacts suggest operational overlap, while reused images may indicate duplicate listings for the same item across different accounts.
- The resulting network provides a visual map of coordinated behaviour, where tightly grouped nodes highlight possible fraud rings.

## Outcome:

- **Several distinct seller communities were detected, each representing possible multi-account networks or duplicate resellers.**
- **These clusters can be explored visually to trace the contact or image links that connect sellers.**

# Graph-Based Clustering-Output

Seller Network — Shared Contacts & Reused Images



## Summary

The graph-based clustering revealed several interconnected seller groups linked through shared contacts and reused product images. Each cluster likely represents coordinated accounts managed by the same individual or team.

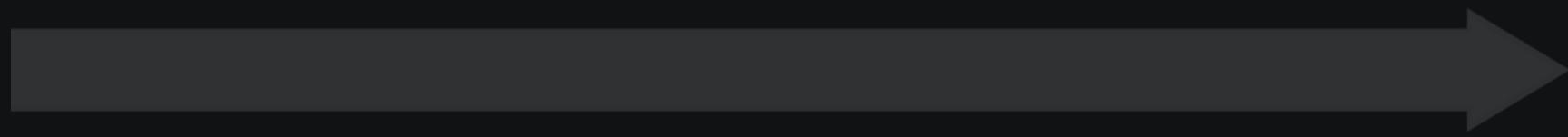
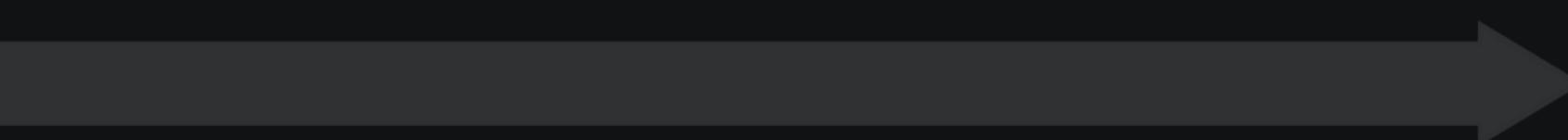
For example, clusters sharing the same email or image URL indicate possible multi-account activity or duplicate listings. For example, Cluster 0 is linked by the contact [mobihub@gmail.com](mailto:mobihub@gmail.com), while others reuse identical product images.

This confirms how graph-based clustering helps uncover hidden seller networks operating under multiple accounts.

# Suspicious Listing Detection using Machine Learning

We built a machine-learning-based detection tool to automatically assess the likelihood of a listing being suspicious.

Using the scraped dataset, the model evaluates indicators such as reused images, shared contacts, suspicious phrases, and pricing anomalies to generate a percentage “risk score” for each listing.



## Model Overview

Algorithm: Isolation Forest (Unsupervised Anomaly Detection)

Framework: scikit-learn

Input Features:

- Price value (price\_num)
- Reused image indicator (reused\_image)
- Shared contact indicator (shared\_contact)
- Price anomaly flag (price\_outlier)
- Seller posting frequency (posts\_per\_day)
- Seller burst activity (seller\_burst\_60m, seller\_burst\_6h)
- High-velocity seller flag (high\_velocity\_seller)
- Duplicate description similarity (dup\_description\_near)
- Output: Suspiciousness Score (0–100%) — representing how different a seller’s behavior is from the norm.

## Results Interpretation

Risk Level	Range (%)	Meaning
Low	0–40	Likely normal seller
Medium	41–70	Needs manual review
High	71–100	Strongly suspicious – reuse/contact overlap detected

# Unsupervised ML for Suspicious Seller Detection

## Isolation Forest: Detecting Suspicious Seller Behavior

### Key Features Used:

- Price anomalies (price\_num, price\_outlier)
- Behavioral metrics (seller\_burst\_60m, posts\_per\_day, high\_velocity\_seller)
- Fraud indicators (reused\_image, shared\_contact, dup\_description\_near)

### Model Overview

- In real marketplaces, ground-truth fraud labels are unavailable.
- Hence, a supervised model (like XGBoost) isn't suitable.
- To overcome this, we used Isolation Forest, an unsupervised anomaly detection algorithm that identifies outliers based on behavioral deviations.
- **Goal: Automatically flag sellers who behave differently from normal patterns.**

### Technical Workflow

- Data preprocessing and feature scaling
- Train Isolation Forest (contamination = 0.1)
- Compute anomaly scores → convert to 0–100 suspiciousness percent
- Assign risk levels → Low, Medium, High

# Unsupervised ML for Suspicious Seller Detection-Ouput

## Isolation Forest: Detecting Suspicious Seller Behavior

Sample of Suspicious Listings:

	seller	price_num	reused_image	shared_contact	suspicious_percent	ml_risk_category
0	aarav78	77016.0	True	False	90.28	High
1	amanigoswami	14235.0	True	False	63.69	Medium
2	amirachadha	122892.0	True	False	83.88	High
3	anahi69	94597.0	True	False	93.19	High
4	anikacherian	141328.0	True	False	82.49	High
5	anyaatwal	120792.0	True	False	84.12	High
6	aradhya39	66725.0	True	False	94.66	High
7	asoni	139285.0	True	False	68.34	Medium
8	azaddas	120250.0	True	False	83.25	High
9	bailaarna	71472.0	True	False	91.10	High

## Observations

The machine learning model assigns a **suspiciousness percentage** to each listing based on indicators like reused images, shared contacts, and pricing anomalies.

Most listings in this sample fall under the **High-risk category (80-95%)**, indicating strong evidence of coordinated or fraudulent behavior.

Only a few listings are rated as **Medium risk**, showing moderate suspicious patterns but not enough overlap for full flagging.

This confirms that the model effectively differentiates between potentially fraudulent sellers and relatively safer ones.

## Outcome:

Overall, this demonstrates the system's capability to automatically identify high-risk sellers and prioritize them for manual verification – reducing human effort and improving fraud detection accuracy.

# OSINT Reconnaissance

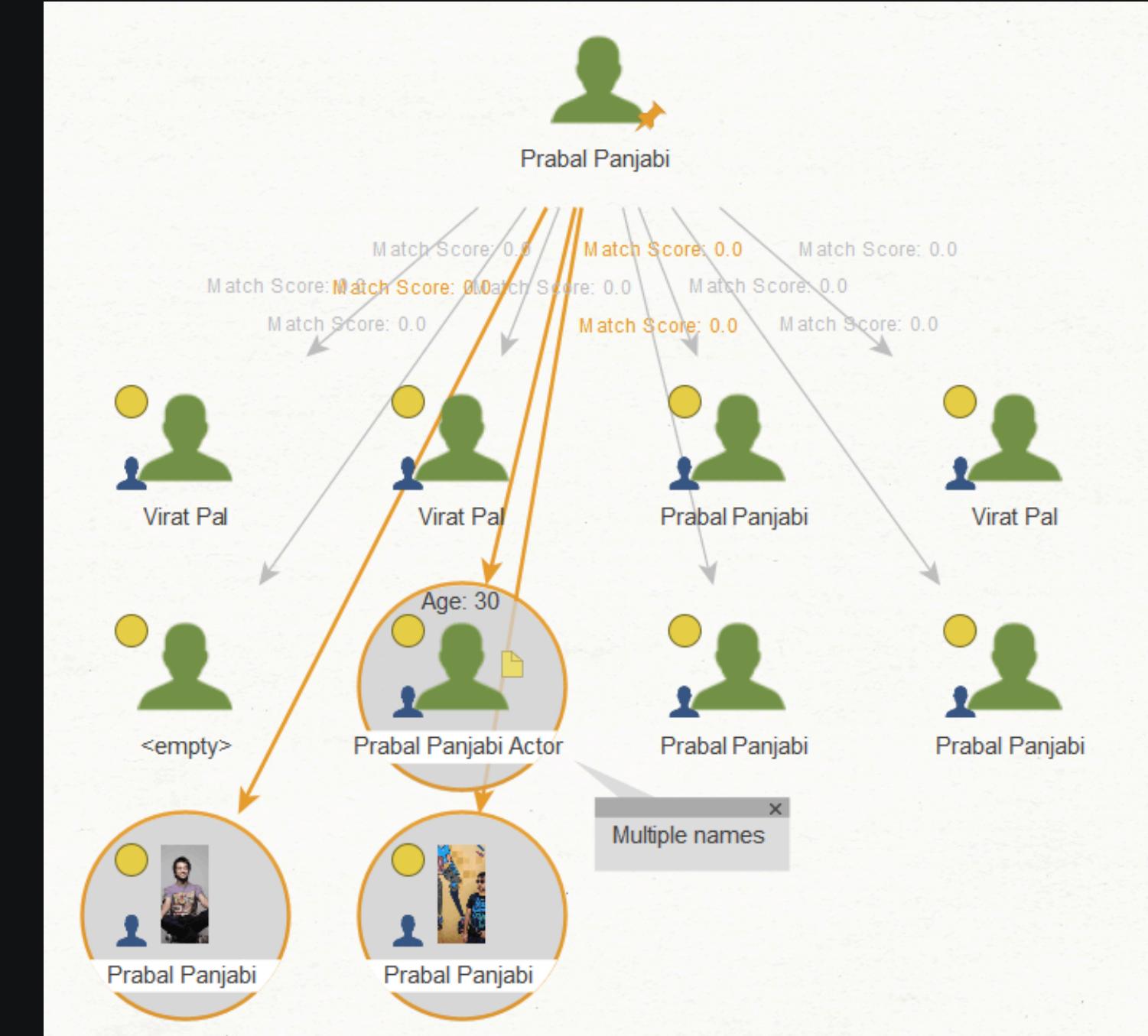
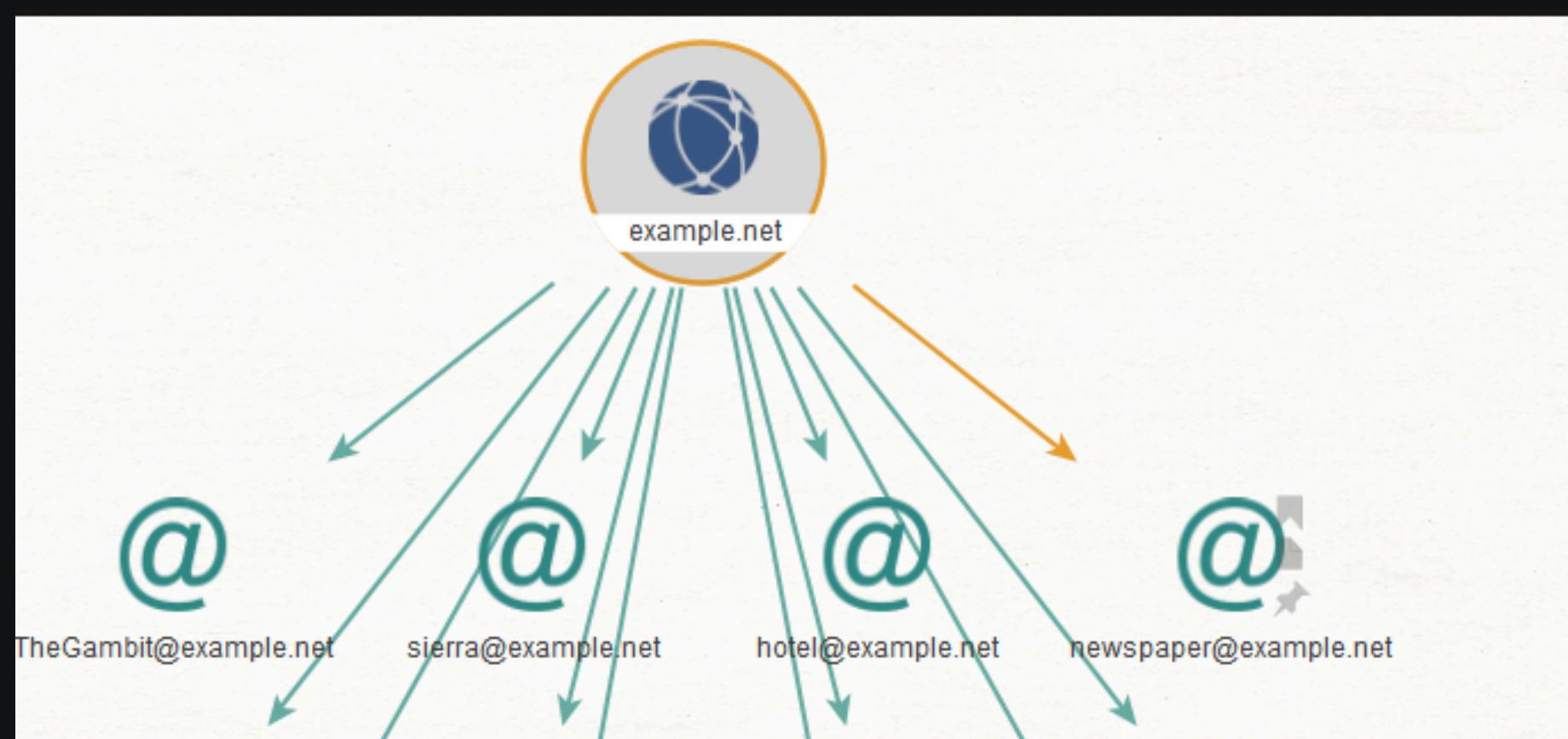
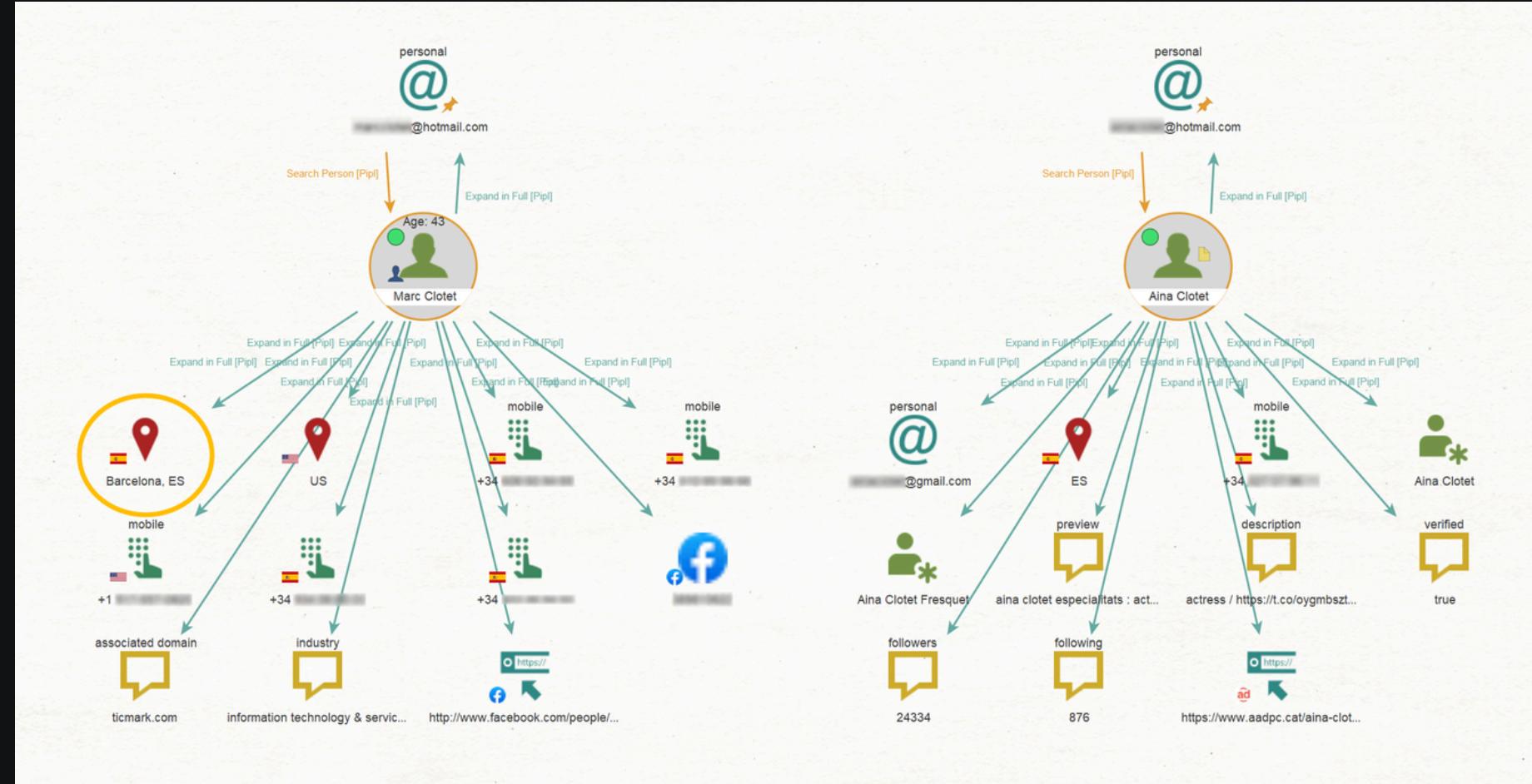
To validate or enrich the suspicious data you scraped (emails, phone numbers, images) using public information sources.

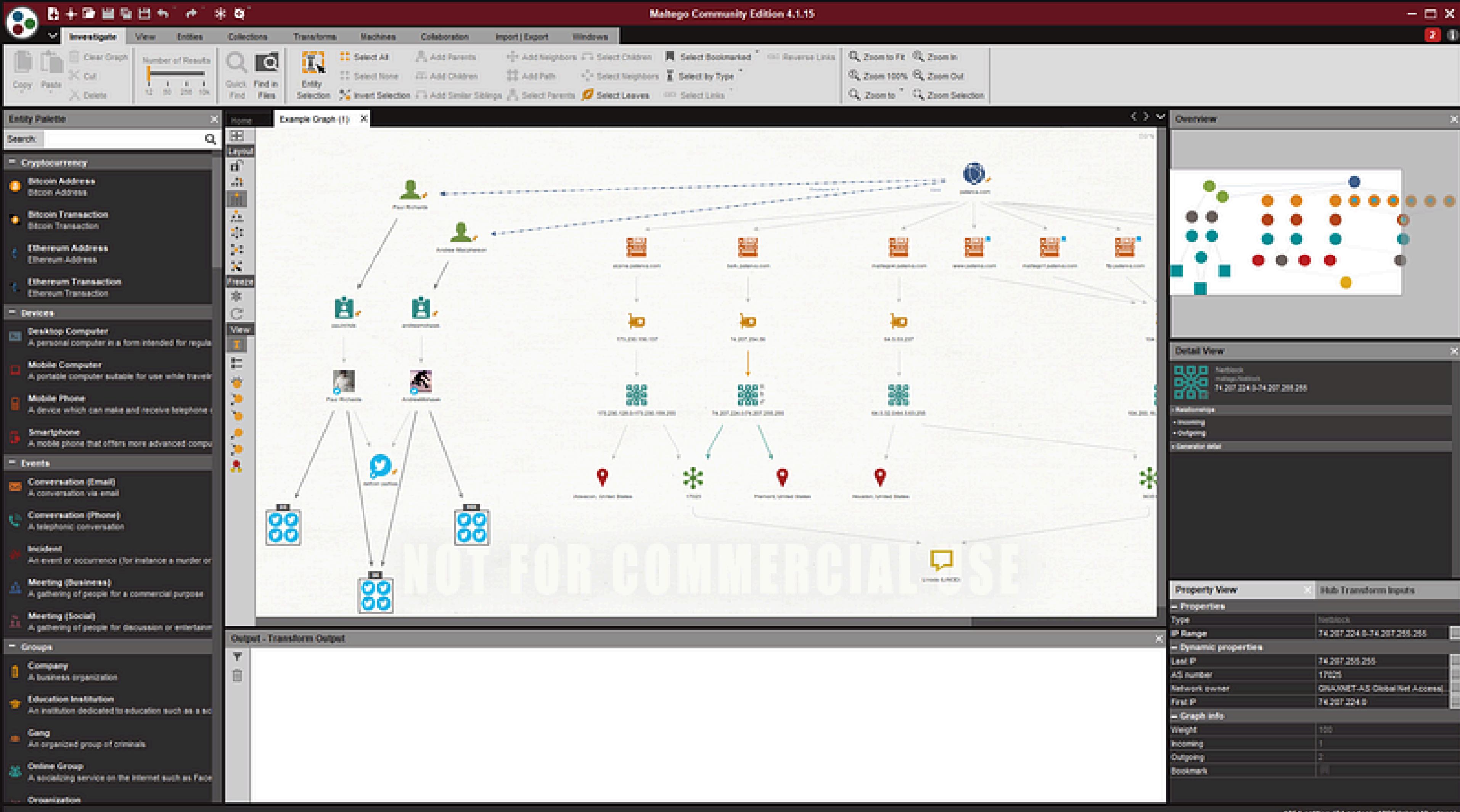
## Maltego

Maltego is an open-source intelligence (OSINT) and link analysis tool used for investigating relationships between digital entities such as emails, domains, and social media accounts. In our project, it helps trace connections between suspicious seller emails and related online footprints, visualizing possible networks of coordinated activity.

## NumVerify

Numverify is a phone number validation and lookup tool that provides details such as carrier, country, and line type. We used it to verify seller contact numbers and identify patterns like reused VoIP or international numbers, which can indicate fake or fraudulent seller accounts.





# Maltego

steps Performed:

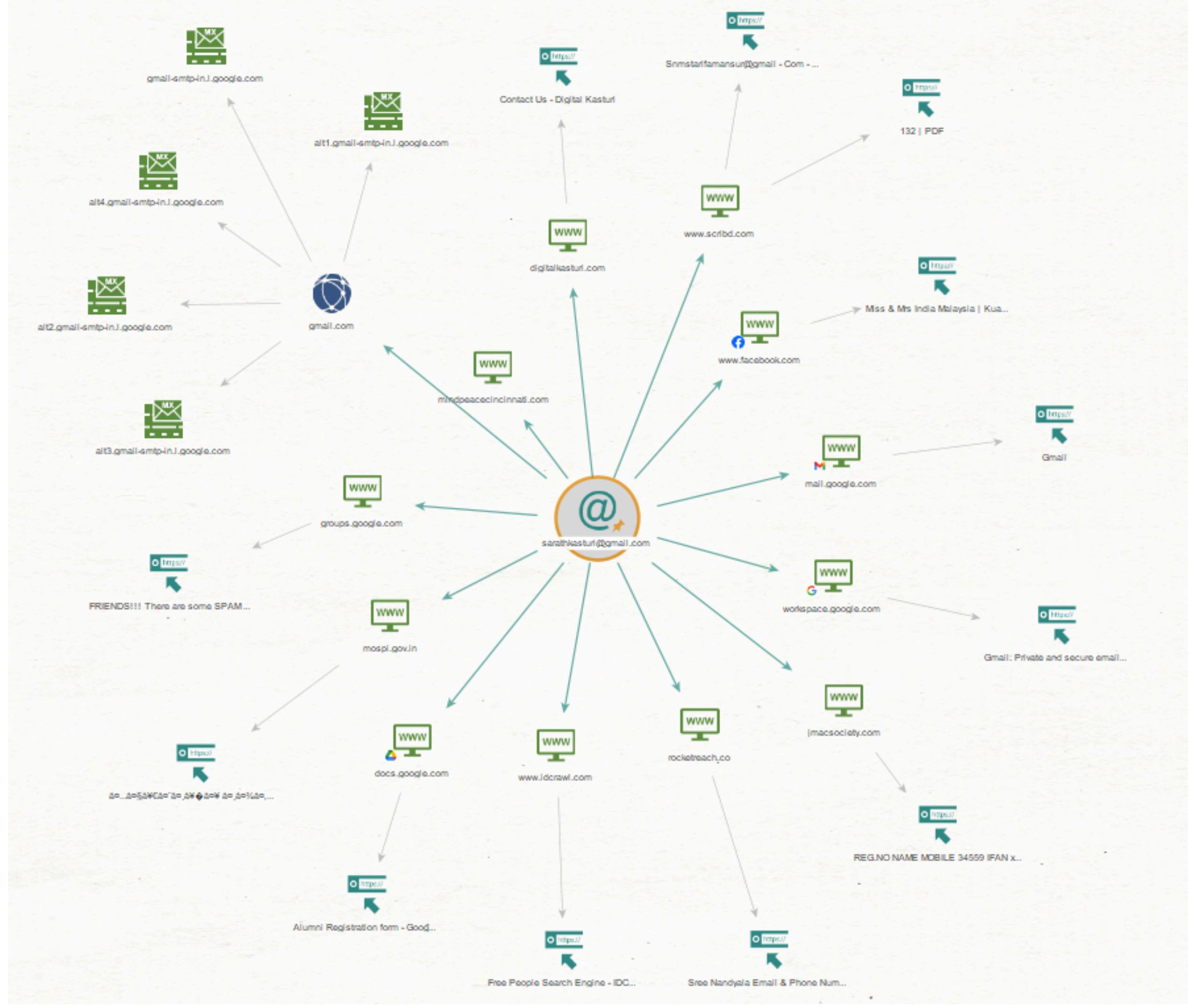
- Added seller email from the scraped dataset.
- Used transforms like “To Domain,” “To MX Record,” “To Entities”.
- Discovered related domains and infrastructure patterns.

Insights Gained:

- Identified sellers using the same domain/email host.
- Found repeated infrastructure footprints (same IP, domain, etc.).
- Helped link multiple listings possibly operated by the same actor.

Use in Investigation:

- ✓ Digital link analysis
- ✓ Network visualization
- ✓ Cross-checking seller identities



# NumVerify

## Steps Performed:

- Queried seller phone numbers using Numverify API.
- Collected metadata: location, carrier, line type (mobile/VoIP).
- Checked for reused or suspicious VoIP numbers.

## Insights Gained:

- Flagged repeated or VoIP-based numbers.
- Detected mismatched country-region codes.
- Strengthened pattern recognition of fake seller networks.

## Use in Investigation:

- ✓ Phone verification
- ✓ Pattern detection
- ✓ Cross-linking reused contacts

 **Valid number!**  
**Country:** India (Republic of)  
**Carrier:** Reliance Jio Infocomm Ltd (RJIL)  
**Location:** Maharashtra

# Conclusion

- Built a full crime-investigation simulation using synthetic marketplace data
- Web scraping automated evidence collection from listings
- Flags & ML detected reused images, contacts, and suspicious listings
- Network visualization (Gephi ) exposed hidden seller clusters & link patterns
- OSINT tools (Maltego, Numverify) enriched investigation with external intelligence
- Demonstrated how data science + OSINT can uncover fraud networks and illicit trade patterns
- All data was synthetic & ethical, proving the workflow is safe and reproducible