

# Neo4j Aura Free 实体去重方案设计

## 方案架构概览

该去重方案包含**离线批处理**和**实时在线**两部分，结合Neo4j Aura Free、Python脚本和n8n工作流来实现。在架构上：

- **Neo4j Aura 图数据库**：存储Person、Company、Report等节点和关系。新增Person节点通过Neo4j HTTP API写入。Aura Free无GDS/APOC插件，所有操作使用原生Cypher和HTTP API。
- **离线批处理去重**：定期（如初始导入或夜间）导出Aura中的Person数据到本地，由Python利用RapidFuzz等工具计算相似度、聚类重复实体，然后通过Neo4j HTTP API批量更新**重复映射**信息回写到图数据库。整个批处理在一小时内完成全图去重。
- **实时在线去重**：在n8n中配置当有新Person通过HTTP API写入Neo4j后，自动触发一个Python脚本流程。该脚本实时查询图中与新节点名称近似的现有Person，利用模糊匹配判断是否重复，并立即写回**规范ID (canon\_id)** 属性或创建**SAME\_AS**关系进行标记。实现新数据“实时”去重映射。
- **前端查询与展示**：数据模型中引入**CanonicalPerson**节点或canon\_id属性以代表实体簇，并保持前端查询体验不变。用户检索时仍可按姓名搜索Person；系统通过Person和CanonicalPerson的联合索引确保重复实体只返回一次结果<sup>1</sup>。

(架构流程：数据源 → n8n HTTP Ingest → Neo4j新增Person → n8n触发Python去重 → Neo4j标注重复 (canon\_id/SAME\_AS) → 前端查询按规范实体聚合显示)

## 离线批处理去重流程

离线批处理用于全量数据定期去重，步骤如下：

1. **数据导出**：通过Neo4j驱动或HTTP API提取所有Person节点的必要字段（至少包含内部ID和name）。例如，执行Cypher查询获取Person清单：

```
MATCH (p:Person)
RETURN id(p) AS person_id, p.name AS name;
```

由于Aura Free无法使用APOC导出，直接用Python调用Neo4j HTTP API或Bolt驱动来获取结果集。得到几万条Person记录后，在本地内存中进行处理。

2. **名称标准化**：对每个Person的name进行清洗和拆分，降低格式差异对匹配的影响。典型处理包括：将名称转小写、去除标点符号、去除多余空格等；按空格/逗号等分隔出name的单词列表并排序。这会得到一个标准化tokens列表（如 name\_tokens），使中间名顺序不同或大小写不同的情况仍能匹配。<sup>1</sup> 离线可直接在Python中处理，实时方案中可在写入前/后触发标准化。

3. **相似度计算 (RapidFuzz)**：利用RapidFuzz库对标准化后的姓名进行模糊匹配打分，找出疑似重复的姓名对。可以采用 `fuzz.token_set_ratio` 或 `fuzz.token_sort_ratio` 作为相似度评分函数，以允许姓名词序调整和大小写差异<sup>2</sup>。具体实现上，为避免对几十万对姓名两两比对的高昂开销，可采取以下策略：
  4. **Blocking策略**：按名称首字母或部分token对数据分块，只在可能重复的子集内计算相似度。例如，将Person按姓氏首字母分组，组内进行模糊匹配；或者根据拆分后的tokens集合交集大小初筛候选。这样减少不可能重复的比较。
  5. **批量匹配**：使用RapidFuzz的批量匹配功能，对于每个姓名，通过 `process.extract()` 在其余姓名列表找出最高相似的若干姓名及分数。针对每个Person保留所有相似度高于阈值（例如85分/100）的匹配。RapidFuzz在C++实现，性能较高，可在阈值剪枝下快速跳过低相似项。
  6. **候选过滤**：如果使用预先算好的 `name_tokens`，也可在Python中快速计算两个姓名tokens集合的交集大小，只有共享至少一个以上token的才送入RapidFuzz精确比对，进一步降低计算量。
7. **形成重复簇 (聚类)**：根据相似度结果，将互相匹配得分超过阈值的Person归为同一重复簇。可建模为无向图：节点代表Person，边代表两者姓名相似度 $\geq$ 阈值的判定。使用NetworkX等工具求该图的连通分量，每个连通分量即为一组重复实体<sup>2</sup>。例如：

```
import networkx as nx
G = nx.Graph()
G.add_nodes_from(person_id_list)
for (id1, id2, score) in similar_pairs: # similar_pairs来自RapidFuzz结果
    G.add_edge(id1, id2)
clusters = list(nx.connected_components(G))
```

得到 `clusters` 列表，其中每个元素是一个集合，包含互为重复的Person节点ID。如果一个簇大小为1，则表示该Person无重复。

8. **分配规范ID (canon\_id)**：为每个重复簇指定一个唯一的**规范标识**。可以为每个簇创建一个新的标识符（例如 `CID12345`，或使用簇中某一代表节点的ID/姓名作为标识）。将此ID作为**canon\_id属性**赋给簇内所有Person节点。如果采用**独立的CanonicalPerson节点模型**，也可为每个簇新建一个 `:CanonicalPerson` 节点，用其 `canon_id` 属性存储唯一标识符，并在Person节点上也写入相同 `canon_id`<sup>3</sup>。这样，每个重复簇由唯一的canon\_id代表。
9. **批量更新回写**：通过Neo4j HTTP API一次性写回所有去重结果，避免频繁网络交互。可构造参数批量设置属性和关系：
10. **创建规范节点 (可选)**：先批量创建所有CanonicalPerson节点及其唯一canon\_id。例如：

```
UNWIND $clusterList AS cluster
CREATE (c:CanonicalPerson {canon_id: cluster.id, name: cluster.name});
```

(其中 `cluster.name` 可选地存一个代表姓名, 如簇中出现频率最高或最完整的姓名)。确保已在Neo4j上建立 `CanonicalPerson(canon_id)` 唯一约束<sup>3</sup>。

11. **设置Person.canon\_id属性**：批量将Person的canon\_id属性更新为对应簇ID。例如使用UNWIND参数列表：

```
UNWIND $personUpdates AS pu
MATCH (p:Person) WHERE id(p) = pu.person_id
SET p.canon_id = pu.canon_id;
```

由于同一簇内Person共享canon\_id, **不能在Person上对canon\_id做唯一约束**, 只在CanonicalPerson上保证唯一性。

12. **建立SAME\_AS关系** (可选)：若不采用canon\_id属性法, 也可使用关系连接重复节点。比如让每个Person连接到对应的CanonicalPerson：

```
MATCH (p:Person),(c:CanonicalPerson)
WHERE p.canon_id = c.canon_id
MERGE (p)-[:SAME_AS]->(c);
```

或直接在重复的Person之间添加无方向的 `SIMILAR` / `SAME_AS` 关系。但推荐使用 **Person→CanonicalPerson** 的映射关系, 方便管理唯一标识。

批处理脚本通过Neo4j官方Python驱动(`neo4j` 库)或HTTP请求执行上述Cypher更新, 实现一次提交多个节点的去重标记<sup>2</sup>。整个离线流程确保在Aura Free的资源限制内高效完成。

## 实时在线去重流程

实时部分通过n8n工作流和Python脚本配合, 使新写入的Person节点立即获得去重标记：

1. **新节点写入触发**：在n8n中配置HTTP接收节点来捕获新Person数据的写入请求（例如来自ETL流程或应用接口）。该节点将Person属性通过Neo4j HTTP API写入Aura（执行类似 `CREATE (p:Person {...})` 的操作）。写入完成后, 紧接着触发后续去重处理。
2. **标准化预处理**：对新节点的name进行同样的标准化处理（小写、去除杂质、拆分token等）。这可以在n8n中完成, 例如使用一个“Function”节点调用JavaScript对名称做简单清洗, 或者在Python脚本里进行。（如果在Neo4j端提前创建了触发器/过程, 也可在节点创建后立即用Cypher对 `p.name_tokens` 等属性赋值, 但Aura Free不支持原生触发器, 这里采用外部处理）。
3. **查询候选相似节点**：新Person名称标准化后, 通过Neo4j查询获取可能重复的现有Person节点列表, 尽量缩小比对范围。由于Aura无全文模糊查询插件, 可利用 **名称token匹配** 作为近似筛选：

```
MATCH (p:Person)
WHERE ANY(token IN $new_name_tokens
```

```

WHERE token IN p.name_tokens)
RETURN id(p) AS person_id, p.name AS name, p.canon_id AS canon_id

```

该查询返回与新节点共享至少一个名字单词的所有Person（即名称有部分重叠）。已对Person.name建立了属性索引或全文索引以加速查询<sup>4</sup>。如果name\_tokens已存为数组属性，Neo4j能够索引数组元素从而较快匹配。必要时可放宽/收紧条件，例如要求至少2个token匹配以提高准确率。

4. **模糊匹配计算**：将上一步返回的候选列表（通常相对较小，如几十个）传递给Python进行相似度计算。n8n可以使用“Execute Command”节点运行Python脚本，或者通过HTTP请求调用一个专门的Python微服务。Python段使用RapidFuzz对**新节点name**与每个候选p.name逐一计算`token_set_ratio`相似度分值，找到最高分和其他高于阈值的匹配：

```

from rapidfuzz import fuzz
best_match = None
best_score = 0
for cand in candidates: # candidates为查询返回的候选人列表
    score = fuzz.token_set_ratio(new_name, cand['name'])
    if score > best_score:
        best_score = score
        best_match = cand

```

设定相似度阈值，如85分以上即认为是重复。如果最佳匹配得分低于阈值，则视为无重复。若有多条候选均超过阈值，也可选取得分最高者作为主要匹配，并考虑其他高分者可能属于同一簇。

5. **确定重复及更新**：根据匹配结果，实时决定如何标记新节点：

6. **找到重复簇**：如果`best_score`高于阈值，说明新节点疑似与现有Person重复：
  - 若匹配候选已有`canon_id`（即已属于某规范簇），则**复用该canon\_id**：通过Neo4j API将新节点的`canon_id`属性设置为相同值<sup>1</sup>。例如：

```

MATCH (p:Person) WHERE id(p) = $new_id
SET p.canon_id = $matched_canon_id;
MERGE (p)-[:SAME_AS]->(:CanonicalPerson {canon_id:
    $matched_canon_id});

```

这样新节点并入已有重复簇。

- 若匹配对象没有`canon_id`（即此前未被标记为重复，意味着这是一个新发现的重复对），则**新建一个规范簇**：生成新的唯一`canon_id`（如新的UUID或组合字符串），并创建对应的CanonicalPerson节点；将该`canon_id`赋给新节点和匹配的旧节点，建立二者指向新CanonicalPerson的SAME\_AS关系。此过程等于动态创建了一个新的重复簇，将原本独立的节点合并进去。
- （如果新节点同时匹配多个候选且分数都高，例如A与B均疑似重复，但A和B本身不同簇，这意味着需要**簇合并**。策略可以是选择其中一个`canon_id`并将另一簇所有成员改写为该ID，从而合并簇。此类情况可留给人工审核或离线批处理二次校正，以免误合并）。

7. **无匹配**：如果没有候选超过阈值，则认为此Person为新唯一实体。可不作任何标记（保持其`canon_id`为空），或为一致性给它分配一个独有的`canon_id`（例如等于其自身ID）并创建一个`CanonicalPerson`节点代表它自己。是否给单例也创建规范ID取决于实现偏好——若前端查询统一走`CanonicalPerson`节点，可以为每个Person都建立对应的`CanonicalPerson`；否则单例可省略`canon_id`。

上述更新通过n8n的HTTP请求节点或Python Neo4j驱动执行。整个实时流程在新节点写入后几秒内完成，确保系统及时标记可能的重复关系。

## 数据建模方案：重复标记与查询

为在不合并节点的前提下标识重复实体，方案采用“规范实体ID + 映射”的建模方式：

- **规范ID ( `canon_id` )**：每组重复Person共享一个字符串ID，用于表示它们属于同一个真实实体。该ID在不同簇间保持唯一。实现上可由系统生成如 `CID1`, `CID2...` 或使用UUID。所有重复成员Person节点的属性 `canon_id` 都设为此值。这样，一个Person如果含有`canon_id`，则说明存在与之重复的其他节点（至少一条）。
- **规范节点 (CanonicalPerson)**：为管理规范ID及保持唯一性，在图中引入新的节点类型 `CanonicalPerson`。每个`CanonicalPerson`节点代表一个实体簇，具有唯一属性 `canon_id` <sup>3</sup>。Person节点通过其`canon_id`与一个`CanonicalPerson`对应，可以选择不存储任何其他信息，或存储该簇的代表姓名、来源等汇总信息。
- **映射关系 (SAME\_AS)**：在模型中增加关系类型 `:SAME_AS`（或 `:SIMILAR`）以连接Person和其对应的`CanonicalPerson`节点。每当一个Person被赋予`canon_id`时，创建 `(p:Person)-[:SAME_AS]->(c:CanonicalPerson)` 关系。这样，重复簇的所有Person都各自指向同一个`CanonicalPerson`，实现聚合表示。也可以从`CanonicalPerson`出发，沿入站`SAME_AS`找到该簇下所有具体Person节点。

**前端查询体验**：通过上述模型，前端查询可以保持原有方式按姓名搜索，但避免返回重复条目：

- 建立一个覆盖Person和`CanonicalPerson`名称的全文索引，使搜索既能匹配具体Person节点名，也能匹配规范节点名 <sup>1</sup>。典型做法是在Neo4j创建联合索引：

```
CALL db.index.fulltext.createNodeIndex('nameFullText',
  ['Person', 'CanonicalPerson'], ['name']);
```

将`CanonicalPerson`的`name`属性设为该簇代表姓名（例如最长的名称或最常用写法）。

- 当用户搜索“张三”时，如果“张三”有多个数据源重复，系统中会有一个`CanonicalPerson[name="张三"]`聚合了重复节点，并出现在全文搜索结果中。同时，由于Person节点也包含“张三”，它们理论上也会匹配；为避免前端看到多个“张三”，可以让查询只返回`CanonicalPerson`节点或对结果去重。一种方式是在搜索查询中筛掉已有`canon_id`的Person，只返回(未归类的Person)和`CanonicalPerson`：

```
CALL db.index.fulltext.queryNodes('nameFullText', '张三') YIELD node
WHERE (node:CanonicalPerson) OR (node:Person AND node.canon_id IS NULL)
RETURN node, labels(node);
```

这样，重复簇以`CanonicalPerson`的形式只出现一次，单例Person正常出现，实现和之前类似的查询体验但无重复显示。前端若需要展开查看具体来源，可根据`CanonicalPerson`→Person的`SAME_AS`关系获取所有对应Person节

点详情。

- 原有依赖Person节点的查询逻辑基本不变。例如，查询某人关系时，依然可以通过Person节点ID查关系；若想涵盖重复，可在查询时先通过`canon_id`找到簇内全部Person再聚合处理。这些都可封装在服务层，保证前端使用感知不到底层去重实现的复杂性。

## 脚本设计与实现

整个方案需若干脚本与自动化流程的支持，均使用Aura兼容的手段（纯Cypher+HTTP）：

- **批量导出脚本（Python）**：连接Neo4j Aura并提取Person数据集。使用Neo4j Python驱动（`neo4j>=5.0` <sup>2</sup>）或requests调用Neo4j HTTP REST接口。示例：

```
from neo4j import GraphDatabase
driver = GraphDatabase.driver("neo4j+s://<Aura-endpoint>",
    auth=("user", "pass"))
with driver.session() as sess:
    result = sess.run("MATCH (p:Person) RETURN id(p) AS id, p.name AS name")
    persons = [record.data() for record in result] # list of dicts
```

数据获取后，进行标准化和RapidFuzz处理，得到重复簇列表及每个Person的归属`canon_id`。

- **重复匹配与聚簇脚本（Python）**：对导出的名单计算模糊相似度并聚类。利用RapidFuzz进行姓名相似度评估 <sup>2</sup> 和NetworkX聚类（或直接用Splink库执行精确的聚类算法）：

```
from rapidfuzz import process, fuzz
names = [p["name_std"] for p in persons] # name_std为标准化后的姓名
# 建立姓名到ID的索引
id_by_name = { p["name_std"]: p["id"] for p in persons }
similar_pairs = []
for p in persons:
    name = p["name_std"]
    # 找到与当前name最相似的前几名候选，阈值筛选
    matches = process.extract(name, names, scorer=fuzz.token_sort_ratio,
        limit=10, score_cutoff=85)
    for match_name, score, _ in matches:
        if score >= 85:
            similar_pairs.append((p["id"], id_by_name[match_name]))
# 基于similar_pairs构建图并找连通分量
G = nx.Graph()
G.add_nodes_from([p["id"] for p in persons])
G.add_edges_from(similar_pairs)
clusters = list(nx.connected_components(G))
```

以上伪码找到所有相似姓名对并归并成重复簇（连通分量）。对于每个簇，生成唯一的`canon_id`标识（如 `CID1` 等）。如果使用Splink（可选），则按配置的字段相似度模型输出簇结果，精度更高但实现复杂度也更高<sup>2</sup>。

- **结果回写脚本（Python）**：将计算得到的簇结果批量更新到Neo4j。可与上一步合并为一个脚本。主要通过Neo4j驱动执行一系列UNWIND批处理语句：

```
# 创建所有CanonicalPerson节点（如采用此模型）
session.run("UNWIND $clusters AS cid CREATE (c:CanonicalPerson {canon_id:
cid})", {"clusters": list(cluster_ids)})
# 更新Person节点的canon_id属性
updates = [{"pid": pid, "cid": cluster_id_map[pid]} for pid in
cluster_id_map]
session.run(

"UNWIND $updates AS u MATCH (p:Person) WHERE id(p)=u.pid SET p.canon_id =
u.cid",
{"updates": updates}
)
```

如果不使用CanonicalPerson节点，也可以省略第一步，仅在Person上设置属性。最后，如需创建关系，可再执行MERGE语句将每个Person连接到其CanonicalPerson。

- **n8n 实时去重 workflow**：在n8n中配置以下节点顺序：
- **Webhook/HTTP接收节点**：接收新Person的数据（JSON/XML），触发流程。
- **Neo4j 创建节点**（HTTP Request节点）：将接收的数据通过Neo4j HTTP API插入，创建新的Person节点。
- **查询候选**（HTTP Request节点）：执行前述Cypher查询，检索与新name有共同token的Person列表。解析Neo4j返回的JSON，提取候选列表。
- **相似度判断**（Function或Execute Command节点）：对候选列表调用Python脚本或JS代码进行RapidFuzz相似度计算，判断最佳匹配及分数。可以将候选列表和新节点name发送到自托管的API服务（如FastAPI应用）进行处理，或者直接在n8n所在主机执行Python脚本<sup>2</sup>。
- **条件判断**（Switch/IF节点）：根据脚本返回结果，分支执行：如果存在重复则进入更新流程，否则结束。
- **更新Neo4j**（HTTP Request节点）：调用Neo4j HTTP API更新新节点的`canon_id`属性，并在需要时创建CanonicalPerson及关系。可直接使用Cypher MERGE/SET完成。例如：

```
// 假设存在匹配簇ID
MATCH (p:Person) WHERE id(p) = $newId
SET p.canon_id = $cid
WITH p
MERGE (c:CanonicalPerson {canon_id: $cid})
MERGE (p)-[:SAME_AS]->(c);
```



如果需要同时更新旧节点（在其之前无`canon_id`的情况），也可一并在Cypher中处理：

```
MATCH (p1:Person) WHERE id(p1)=$newId
MATCH (p2:Person) WHERE id(p2)=$matchedOldId
CREATE (c:CanonicalPerson {canon_id:$newCid})
SET p1.canon_id=$newCid, p2.canon_id=$newCid
MERGE (p1)-[:SAME_AS]->(c)
MERGE (p2)-[:SAME_AS]->(c);
```

通过n8n这样编排，实现新节点创建→查重→标记的一条龙自动处理。

以上脚本和流程均只使用Neo4j的Cypher查询与HTTP接口，**未依赖任何Aura不支持的插件**（未使用APOC存储过程或GDS算法），完全兼容Aura Free环境要求。开发时需确保Python环境安装了所需库（如RapidFuzz、Neo4j驱动、NetworkX等<sup>2</sup>）。

## 可选增强：Streamlit前端审核

为了提高去重准确性，可以添加一个**Streamlit轻量前端**供人工审核和override结果（选配）：

- **聚簇结果检查**：Streamlit应用连接Neo4j数据库（使用Neo4j Python驱动）读取所有已标记的重复簇（例如查询所有CanonicalPerson及关联的Person）。将每个簇的成员姓名、来源等显示在网页表格中。人工可以浏览哪些Person被归为同一簇。
- **人工校正**：对于怀疑错误的归并，允许用户手动调整。例如提供“拆分”按钮，将选定Person从当前簇移除（即重新赋予新`canon_id`或清除其`canon_id`）；“合并”按钮，可选中两个规范簇合并为一（为它们指定相同`canon_id`）。这些操作背后通过调用Neo4j的HTTP API执行更新。
- **参数调整和重跑**：在Streamlit界面上提供阈值调整选项，允许用户修改RapidFuzz相似度阈值或规则，然后触发重新计算（调用离线脚本重新跑一遍去重逻辑），以观察不同参数对结果的影响，辅助找到最佳阈值设置。
- **实时监控**：Streamlit页面也可显示最新写入的Person记录及其自动匹配结果，方便人工及时审阅。如果发现误判，可立即通过界面纠正（例如将错误合并的节点`canon_id`清空或改为正确簇）。

通过上述人工审核界面，系统管理员能对自动去重结果进行把关和调整，从而不断完善知识图谱的实体质量。

**总结**：这套方案在Neo4j Aura Free环境下实现了实体去重：利用Python+RapidFuzz离线计算全图重复簇并批量标记<sup>5</sup>；通过n8n流水线结合Cypher查询和Python脚本，实现新数据的实时去重映射；采用`canon_id`属性和CanonicalPerson节点建模重复关系，在不合并节点的情况下将重复实体关联起来，并保持前端查询体验一致<sup>1</sup>。所有步骤仅借助Cypher和HTTP API完成，符合Aura限制。必要时增加简易前端供人工复核，进一步提高方案可靠性和可维护性。这样即可获得一套**可直接部署**的Neo4j知识图谱实体去重解决方案。<sup>4</sup> <sup>1</sup>

---

<sup>1</sup> <sup>2</sup> <sup>3</sup> <sup>4</sup> <sup>5</sup> Dedup Pipeline Full Code.docx

<https://drive.google.com/file/d/1XF2dHqf6v1NSPNF3XAmbIIYBnNS7iCXj>