

Compte rendu — Workflow n8n : Recherche de stages à l'étranger

Leo Torres — DO3

19 janvier 2026

Table des matières

1	Contexte et objectifs	2
1.1	Objectif du projet	2
1.2	Périmètre et hypothèses	2
2	Architecture du workflow n8n	2
2.1	Vue d'ensemble du workflow	2
2.2	Description des étapes	2
2.2.1	Étape 1 — Génération des combinaisons ville × technologie	2
2.2.2	Étape 2 — Recherche web (Tavily)	2
2.2.3	Étape 3 — Enrichissement scoring via MCP	2
2.2.4	Étape 4 — Extraction d'informations via LLM	3
2.2.5	Étape 5 — Génération des résumés via LLM	3
2.2.6	Étape 6 — Filtrage personnalisé	3
2.2.7	Étape 7 — Export des résultats	3
2.2.8	Étape 8 — Synthèse finale	3
2.2.9	Étape 9 — Notification	3
3	Exemples de résultats	3
3.1	Offres extraites (5 à 10 exemples)	3
3.2	Fichier CSV exporté (ou capture)	4
3.3	Workflow n8n (JSON) dans le repo	4
4	Bonus — Méthode de scoring	5
4.1	Méthode utilisée	5
4.2	Choix de scoring et justification	5
4.3	Évaluation : est-ce que ça fonctionne comme attendu ?	5
5	Analyse critique	5
5.1	Difficultés rencontrées	5
5.2	Critères de filtrage choisis et justification	5
5.3	Apports du workflow et axes d'amélioration	5
6	Conclusion	6
A	Annexes	6
A.1	Prompts LLM (extraits)	6

1 Contexte et objectifs

1.1 Objectif du projet

Le workflow automatise la recherche d'offres de stage a l'étranger a partir d'une matrice *ville × technologie*. Les resultats sont enrichis par un score de ville (MCP), normalises par LLM, dedoublonnes, puis exportes vers Google Sheets et resumes dans une synthese envoyee sur Discord.

1.2 Périmètre et hypothèses

- Villes : Berlin, Stockholm.
- Technologies : crypto, devops, cybersecurity, iot (8 combinaisons).
- Source principale : recherche web via Tavily (1 resultat par requete).
- LLM : Google Gemini (gemma-3-4b-it pour extraction, gemma-3-27b-it pour resume, gemini-robotics-er-1.5-preview pour synthese) (à changé en fonction de la dispo des quotas gratuit).
- Limites : dependance aux API, champs souvent absents (salaire, remote), redondances d'URLs.

2 Architecture du workflow n8n

2.1 Vue d'ensemble du workflow

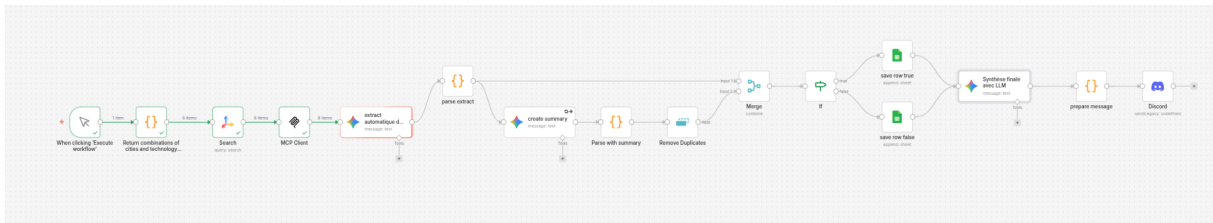


FIGURE 1 – Workflow n8n complet.

2.2 Description des étapes

2.2.1 Étape 1 — Génération des combinaisons ville × technologie

Un node Code genere les couples ville x technologie a partir de deux tableaux statiques.

2.2.2 Étape 2 — Recherche web (Tavily)

Chaque couple declenche une recherche Tavily avec la requete type `company <tech> <city> internship`. La recherche est en `basic`, avec `max_results=2`.

2.2.3 Étape 3 — Enrichissement scoring via MCP

Le MCP Client appelle l'outil `find_city_score` avec le nom de ville. Si la ville est connue, un score est renvoye et ajoute aux donnees d'entree ; sinon le score est absent et le MCP indique `found=false`.

2.2.4 Étape 4 — Extraction d'informations via LLM

Le LLM retourne un JSON strict, sans texte additionnel, conforme a un schema fixe (company, job_title, city, country, remote_policy, contract_type, salary, currency, duration, application_deadline, skills, languages, source_url, source_title, city_score). Les valeurs manquantes sont forcees a null.

2.2.5 Étape 5 — Génération des résumés via LLM

Un second LLM produit un resume FR de 2 a 3 phrases par offre. Les sorties non-JSON sont nettoyees et parsees dans un node Code.

2.2.6 Étape 6 — Filtrage personnalisé

Le filtrage est minimal : dedoublonnage par `source_url`, puis routage selon la presence d'un salaire pour separer deux onglets de stockage.

2.2.7 Étape 7 — Export des résultats

Les offres sont ecrites dans Google Sheets (equivalent CSV), dans deux onglets en fonction du filtre (pour les offres ok selon le filtre et pour les autres). Colonnes : company, job_title, contract_type, city, country, salary, currency, duration, application_deadline, skills, languages, source_url, source_title, city_score, summary, remote_policy.

2.2.8 Étape 8 — Synthèse finale

Un LLM genere une synthese lisible avec un titre, le nombre d'offres, un top 3 argumente et des statistiques globales (villes, contrats, remote, salaires).

2.2.9 Étape 9 — Notification

La synthese est envoyee sur Discord via webhook. Un node Code tronque le message a 1900 caracteres pour respecter la limite.

3 Exemples de résultats

3.1 Offres extraites (5 à 10 exemples)

Exemple de 5 a 10 offres extraites (toutes les infos dans le JSON) :

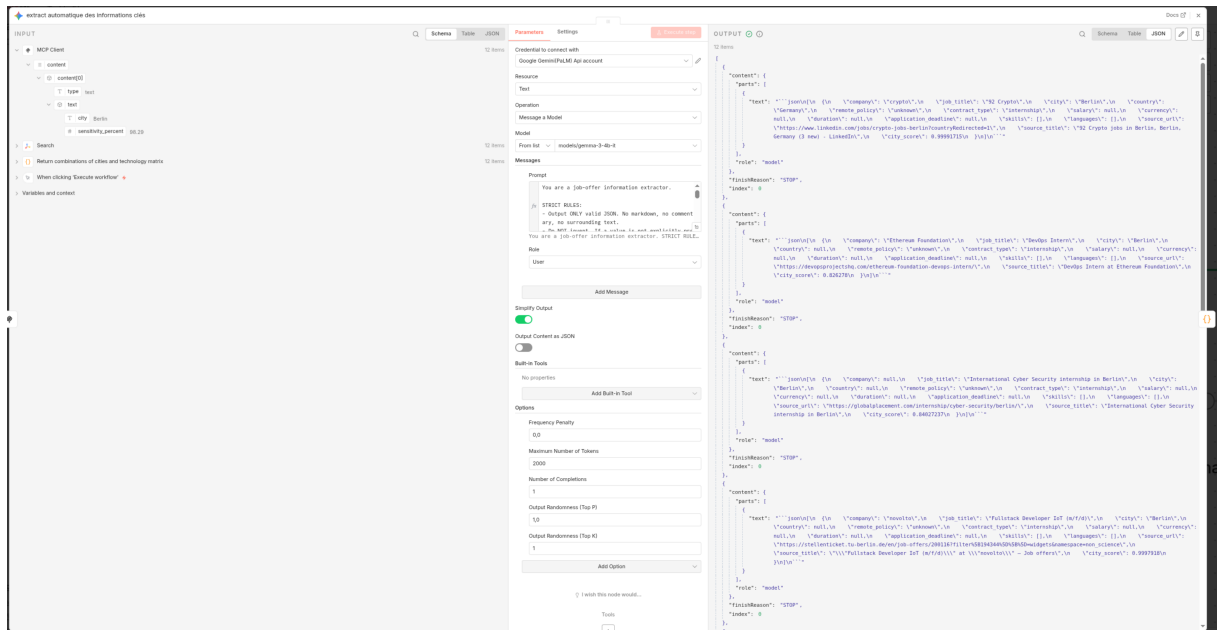


FIGURE 2 – Exemple de 5-10 offres extraites (JSON).

3.2 Fichier CSV exporté (ou capture)

company	job_title	city	country	remote_policy	contract_type	salary	currency	duration	application_deadline	skills	languages	source_url	source_title	city_score	summary
Centiglobe	Summer intern	Stockholm		unknown	internship	Competitive				[Blockchain, crypto]		https://hellub.se	Summer intern	0.78944516	Centiglobe propose un stage d'été en Blockchain et cryptomonnaie à Stockholm. Il s'agit d'un c
crypto	Web3 internship	Barcelona	Spain	unknown	internship	\$30k/year to \$120k/year						https://hellub.se	Web3 Internship	0.869603	Crypto propose des stages Web3 à Barcelone, en Espagne. Il s'agit d'un contrat de stage avec

FIGURE 3 – Aperçu de l'export CSV / Google Sheets.

3.3 Workflow n8n (JSON) dans le repo

```

1 repo/
2   workflow-final.json

```

4 Bonus — Méthode de scoring

4.1 Méthode utilisée

- Source : table interne `CITY_SENSITIVITY` avec des valeurs S_j (en %) par ville.
- Calcul : transformation *inverse min-max* pour obtenir un score 0–100 (plus lisible).
- Acces via MCP : outil `find_city_score` expose par le serveur MCP.
- Correspondance : recherche insensible a la casse (ex. *berlin* = *Berlin*).
- Valeurs manquantes : si la ville n'est pas dans la table, l'outil renvoie une erreur (pas de score).

Calcul utilise :

$$\text{city_score} = \left(1 - \frac{S_j - \min(S_j)}{\max(S_j) - \min(S_j)} \right) \times 100$$

Le score est arrondi a 2 decimales. Un S_j faible (forte concordance) donne un score plus eleve, tandis qu'un S_j eleve est penalise.

Le calcul est implemente directement dans le serveur MCP.

4.2 Choix de scoring et justification

- Pertinence : indicateur externe pour comparer rapidement l'attractivite des villes.
- Simplicité : table pre-calculée, sans besoin de recalcul complexe.
- Robustesse : le score est stable et ne depend pas des offres collecte es.

4.3 Évaluation : est-ce que ça fonctionne comme attendu ?

Pour les villes presentes dans la table, le score est coherent et stable. Les limites viennent des villes absentes ou des libelles atypiques (ex. quartiers, regions). Les aliases ameliorent le taux de correspondance, mais un enrichissement geographique serait utile.

5 Analyse critique

5.1 Difficultés rencontrées

- Choix d'un LLM gratuit : trouver un modele avec un plan gratuit, assez rapide et suffisamment performant pour des extractions fiables.
- Parsing et fusion : difficultes a parser les sorties LLM puis a merger les flux pour n'obtenir qu'un seul objet exploitable.
- Webhook Discord : configuration et tests pour respecter les contraintes de taille des messages.
- Google Sheets : configuration Google Cloud API (OAuth) plus longue que prevu avant de pouvoir ecrire dans un tableur.

5.2 Critères de filtrage choisis et justification

Le filtrage est volontairement leger pour maximiser la couverture : dedoublonnage par URL et separation selon `salary` present ou non. Cela permet d'identifier rapidement les offres les plus informatives sans exclure trop de candidats.

5.3 Apports du workflow et axes d'amélioration

- Gain de temps important sur la collecte et la mise en forme des offres.
- Ameliorations : augmenter `max_results`, ajouter des sources, filtrer par `city_score`, enrichir la deduplication, meilleure gestion d'erreurs et retries.

6 Conclusion

Le workflow automatise efficacement la recherche et la synthèse d'offres de stage. L'extraction structurée et l'export permettent une analyse rapide, tandis que le scoring MCP apporte un signal supplémentaire. Les principaux axes d'amélioration concernent la couverture des sources et la robustesse face aux données incomplètes.

A Annexes

A.1 Prompts LLM (extraits)

```
1 Prompt d'extraction (Gemini) :
2 You are a job-offer information extractor.
3
4 STRICT RULES:
5 - Output ONLY valid JSON. No markdown, no commentary, no surrounding
  text.
6 - Do NOT invent. If a value is not explicitly present, use null.
7 - Empty lists must be [].
8 - Extract ONE output object per offer in the input array, preserving the
  same order and length.
9
10 OUTPUT FORMAT:
11 Return a JSON array of objects. Each object must follow EXACTLY this
  schema:
12
13 {
14   "company": string|null,
15   "job_title": string|null,
16   "city": string|null,
17   "country": string|null,
18   "remote_policy": "onsite"|"hybrid"|"remote"|"unknown"|null,
19   "contract_type": "internship"|"apprenticeship"|"full_time"|"
    part_time"|"unknown"|null,
20   "salary": string|null,
21   "currency": string|null,
22   "duration": string|null,
23   "application_deadline": string|null,
24   "skills": string[],
25   "languages": string[],
26   "source_url": string|null,
27   "source_title": string|null,
28   "city_score": number|null
29 }
30
31 MAPPING RULES:
32 - source_url = offer.url if present, else null
33 - source_title = offer.title if present, else null
34 - city_score = offer.score if present, else use input.search_score if
  provided, else null
35 - remote_policy:
36   - "remote" if clearly remote
37   - "hybrid" if clearly hybrid
38   - "onsite" if clearly onsite/on-site
39   - otherwise "unknown"
40 - contract_type:
41   - "internship" if stage/internship/intern
```

```

42 - "apprenticeship" if alternance/apprenticeship
43 - otherwise "unknown"
44 - salary: keep original text (e.g. "800 EUR/month", "$22k-$62k")
45 - currency: "EUR", "USD", "GBP" if clearly identifiable, else null
46
47 INPUT:
48 You will receive a JSON object containing offers in an array field. Use:
49 - offers = input.result OR input.results (whichever exists)
50
51 Input JSON :
52 {
53     "city": string,
54     "sensitivity_percent": number,
55 }
56
57 -----
58
59 Prompt de resume (Gemini) :
60 Generate a short readable summary for each extracted offer.
61
62 STRICT RULES:
63 - Output ONLY valid JSON. No markdown, no commentary.
64 - Return a JSON array with the same length/order as the input array.
65 - Do not include any other keys.
66
67 OUTPUT SCHEMA (per item):
68 {
69     "source_url": string|null,
70     "summary": string
71 }
72
73 SUMMARY RULES:
74 - 2 to 3 sentences maximum.
75 - Include: company (if known), role, location, remote policy, contract
76     type, and salary if present.
77 - If a field is missing, omit it (do not guess).
78 - Write the summary in French.
79
80 INPUT JSON (array of extracted offers):
81 [
82     {
83         "company": "crypto",
84         "job_title": "Director of Venture Studio",
85         "city": "Berlin",
86         "country": "Germany",
87         "remote_policy": "unknown",
88         "contract_type": "internship",
89         "salary": null,
90         "currency": null,
91         "duration": null,
92         "application_deadline": null,
93         "skills": [],
94         "languages": [],
95         "source_url": "https://www.linkedin.com/jobs/crypto-jobs-berlin?countryRedirected=1",
96         "source_title": "95 Crypto jobs in Berlin, Berlin, Germany (1 new)",
97         "city_score": 0.99972624

```

```

97     },
98     ...
99 ]
100
101 OUTPUT SCHEMA (per item):
102 {
103     "source_url": string|null,
104     "summary": string
105 }
106
107 -----
108
109 Prompt de synthese finale (Gemini) :
110 You are a report writer for internship search results.
111
112 STRICT RULES:
113 - Output ONLY plain text. No JSON, no markdown, no commentary.
114 - Write in French.
115 - Keep the message under 1900 characters.
116 - Do not invent numbers or facts.
117
118 INPUT:
119 You will receive a JSON array of extracted offers. Each item may include
120 :
121 company, job_title, city, country, remote_policy, contract_type, salary,
122     currency, duration, application_deadline, skills, languages,
123     source_url, source_title, city_score, summary.
124
125 OUTPUT FORMAT:
126 - Title on the first line.
127 - Then 3 to 5 short paragraphs:
128     1) Nombre total d'offres valides.
129     2) Top 3 des offres les plus interessantes avec 1 phrase de
130         justification chacune (score ville, salaire, remote, contrat,
131         pertinence tech).
132     3) Statistiques globales (villes les plus presentes, types de contrats
133         , remote, salaires disponibles).
134
135 RULES:
136 - If fewer than 3 offers, list all.
137 - If a field is missing, skip it.
138 - Use source_url only if needed to disambiguate.

```