

Automatic data extraction from EFSA Opinions

Diana Varšíková, Ecomole s.r.o

August 27, 2024

1 Introduction

We received csv file with 204 question numbers/opinions to include in the database. Below is an example of the data that are to be extracted:

Administrative Data of the Opinion: Applicant company, their country of origin, DOI of the opinion, plus 10 other fields.

General Information: Trade names, common names, form of the food (whole food, extract, etc.).

Identity: This depends on the category. For example, for animals: genus, species, subspecies, part used; for chemicals: common name, IUPAC name, CAS number, molecular formula, etc.

Production Process: List of main production steps.

Main Composition of the Novel Food: Carbohydrates, proteins, fats, minerals, water, vitamins.

Proposed Uses and Use Levels: Proposed uses (whole foods, supplements, etc.), target population (general, infants, children, etc.), food category (FoodEx / FAIM).

Nutritional Information: Nutritionally disadvantageous / nutritionally advantageous components.

Availability of ADME (Absorption, Distribution, Metabolism, and Excretion) Studies.

Availability of Toxicological Studies and Their Outcome.

Allergenicity Assessment: Unlikely, low, possible, certainty etc.

The primary objective was to automate the extraction of 'non-expert' information—data that can be easily understood and assessed for accuracy by individuals without a biology degree. This approach allows experts to focus on more complex and technical information. The goal was to automatically extract the majority of the administrative data and some elements of the general information. This report will detail which data were successfully extracted automatically and the logic and methods used in the extraction process.

2 Data Summary

We received a CSV file from EFSA containing a list of 207 question numbers to be included in the database. The types of articles associated with these question numbers are as follows:

Type Of Article	Count
Opinion	198
Guidance	4
Public consultation	3
Data gaps statement	1

This distribution indicates that 198 opinions are to be extracted, as it was agreed not to extract data from the other types of articles. Out of the 198 opinions, EFSA provided Ecomole with JATS files for 194 of them. However, there are variations in the completeness of these files:

JATS availability	count
full JATS file	123
front JATS file	71
no JATS file	4

This summary highlights that 71 of the JATS files contain only the FRONT section, which includes administrative information but lacks the content part necessary for a full extraction.

3 Opinions analysis


By looking at the opinions, the structure of the opinions can be divided into five groups.

Group A

This is the most common group.

<p>SCIENTIFIC OPINION</p> <p>ADOPTED: 13 December 2016 doi: 10.2903/j.efsa.2017.4682</p> <p>Scientific Opinion on taxifolin-rich extract from Dahurian Larch (<i>Larix gmelinii</i>)</p> <p>EFSA Panel on Dietetic Products, Nutrition and Allergies (NDA), Dominique Turck, Jean-Louis Brosson, Barbara Burlingame, Tara Dean, Susan Fairweather-Tait, Marina Heinrich, Gábor István Hirsch-Estlin, Ugo Marquardt, Harry J. Mulder, Androniki Naska, Monika Neuhäuser-Berthold, Grzyzyna Nowicka, Kristina Penttinen, Yolanda Sanz, Alfonso Sanz, Anders Skjold, Martin Stern, Daniel Tsiang, Marco Vioroli, Peter Willatts, Karl Heinz Engel, Rosangela Marchelli, Arnette Pölzig, Morten Poulsen, Josef Schlatter, Wolfgang Gebmann and Henk van Loveren</p> <p>Abstract</p> <p>Following a request from the European Commission, the EFSA Panel on Dietetic Products, Nutrition and Allergies (NDA) was asked to carry out the additional assessment for taxifolin-rich extract from Dahurian Larch as a food ingredient in the context of Regulation (EC) No 2581/07. The novel food (NF) is a taxifolin-rich water-ethanol extract from the wood of the Dahurian Larch and contains a minimum of 80% taxifolin. The Panel considers that the taxifolin-rich extract is sufficiently characterised and that its compositional data and specifications do not raise safety concerns. The NF is intended to be added to non-alcoholic beverages, to yogurt and to chocolate confectionery. The Panel considers that the data on genotoxicity do not raise concern. In a subchronic rat study performed in accordance with OECD standards, the highest dose tested (i.e. 1,500 mg/kg bw) was considered to be the NOEL. The margin of exposure (MOE) of the combined intake (158 mg) from the intended food uses (including 100 mg from food supplement) would result in about 600 for an adult weighing 70 kg. For adolescents, taking into account a default body weight of 45 kg, the MOE of the combined intake (146 mg) would be about 600. In the absence of a high percentile intake estimate for children between 9 and 14 years of age, the Panel considers the P97.5 intake estimate from the intended food uses (except from food supplements) for children between 10 and 17 years, i.e. 46 mg/day, taking into account a default body weight of 29.4 kg (PS body weight for children aged 10–14 years as suggested by EFSA Scientific Committee (2012)), the resulting MOE would be about 960.</p> <p>© 2017 European Food Safety Authority. EFSA Journal published by John Wiley and Sons Ltd on behalf of European Food Safety Authority.</p> <p>Keywords: taxifolin, (2R,3R) trans-dihydroquercetin, 2,3-dihydroquercetin, Dahurian Larch, novel food, ingredient</p> <p>Requester: European Commission following an application by Ametis SC</p> <p>Question number: EFSA-Q-2012-0961</p> <p>Correspondence: nda@efsa.europa.eu</p>	<p>Safety of taxifolin</p> <p>Table of contents</p> <p>Abstract..... 1</p> <p>Summary..... 2</p> <p>1. Introduction..... 3</p> <p>1.1. Background and Terms of Reference as provided by the European Commission..... 3</p> <p>1.2. Data and Methodologies..... 5</p> <p>2.1. Data..... 6</p> <p>2.2. Methodologies..... 6</p> <p>3. Assessment..... 6</p> <p>3.1. Specification of the NF..... 6</p> <p>3.2. Effect of the production process applied to the NF..... 8</p> <p>3.3. History of the exposure used in the course of the NF..... 8</p> <p>3.4. Anticipated intake/level of use of the NF..... 8</p> <p>3.5. Information from previous exposure to the NF..... 10</p> <p>3.6. Nutritional information on the NF..... 10</p> <p>3.7. Microbiological information on the NF..... 10</p> <p>3.8. Toxicological information on the NF..... 10</p> <p>3.8.1. Genotoxicity..... 10</p> <p>3.8.2. Absorption, distribution, metabolism, excretion..... 11</p> <p>3.8.3. Acute and subacute toxicity..... 12</p> <p>3.8.4. Subchronic toxicity..... 12</p> <p>3.8.5. Developmental toxicity..... 13</p> <p>3.8.6. Reproductive toxicity..... 14</p> <p>3.8.7. Allergenicity..... 14</p> <p>4. Discussion..... 14</p> <p>5. Conclusions..... 14</p> <p>Documentation provided to EFSA..... 14</p> <p>References..... 15</p> <p>Abbreviations..... 16</p>
---	--

The image shows the front cover of a report from the European Food Safety Authority (EFSA). At the top left is the EFSA logo, which consists of the word 'efsa' in a stylized blue font with a blue square icon to its right, and the full name 'European Food Safety Authority' in a smaller blue font below it. To the right of the logo, the text 'The EFSA Journal (2008) 1(2)' is printed in a small, black, sans-serif font. The main title of the report, 'Safety of Synthetic Lysine¹', is centered at the top in a large, bold, black, sans-serif font. Below the title, the subtitle 'Panel on Safety on Panel on Scientific Panel on Dietetic Products, Nutrition and Allergies' is centered in a smaller, bold, black, sans-serif font. Underneath the subtitle, the text '(Question No EFSA-Q-2007-119)' is centered in a smaller, regular, black, sans-serif font. The next line of text, 'Adopted on 10 April 2008 by written procedure', is centered in a bold, black, sans-serif font. Below this is a horizontal line. Under the line, the text 'PANEL MEMBERS' is centered in a bold, black, sans-serif font. This is followed by a list of names: 'Jean-Louis Brethes, Albert Flynn, Martin Hejblum, Karin Hubshä, Hanna Korhonen, Pagona Lagou, Martinus Leuz, Rosangela Mucchelli, Ambrose Martin, Bessy Mvondo, Andreus Palm, Hildegard Rydzynski, Sappho Salmassi, John Sato Strate, Stephan Strobel, Inge Thomsen, Heek van der, Hendrik van Loveren, and Hans Verhagen.' Below the list of names is another horizontal line. Under this line, the text 'SUMMARY' is centered in a bold, black, sans-serif font. The main body of the text is in a regular, black, sans-serif font. It begins with 'Following a request from European Commission, the Panel on Dietetic Products, Nutrition and Allergies was asked to deliver a scientific opinion on the safety of synthetic lysine for use as and proposed as a lysine.' The text continues: 'The applicant proposes to use synthetic lysine both as a food supplement and as a food ingredient. The novel food ingredient consists of synthetic (crystalline) lysine that is marketed in three different formulations. These are lysineose 10%, lysineose 10 cold water dispersible (CWD) and lysineose dispersion 20%.' The text then continues: 'Synthetic lysine is suggested by the applicant to be used in food supplements at levels of 8 to 15 mg/kg, in beverages and dairy products at levels of up to 2.2 mg/100 g, in bread/cereals up to 4 mg/100 g, in special bars up to 8 mg/100 g, in fats and dressings up to 4 mg/100 g and in dietary foods in general by mass at levels in accordance with the particular requirements of the person for whom the product is intended.' The text concludes with: 'The applicant provides an intake estimate of lysineose based on three sources including 1) normal dietary intake from food, 2) intake from dietary supplements and 3) intake from proposed fortified food products. The Panel notes that an additional source is 4) use as a food colorant.'

	<i>Salmonella</i> <i>Salmonella</i> <i>Salmonella</i>
<hr/>	
TABLE OF CONTENTS	
Panel Members	1
Summary	1
Table of Contents	1
Background as provided by the commission	4
Terms of reference as provided by the commission	6
Acknowledgements	7
Assessment	8
I. Specification of the novel food (NF)	8
II. Effect of the production process applied to the NF	10
III. History of the operation used as the source of the NF	11
IV. Anticipated industrial use of the NF	11
V. Information from previous human exposure to the NF or its source	14
VI. Nutritional information on the NF	14
VII. Microbiological information on the NF	15
VIII. Toxicological information on the NF	15
Discussion	19
Conclusions and Recommendations	22
Documentation provided to EFSA	22
References	22

Group E

4.1 DOI

The DOI can be found in article-meta - article-id - pub-id-type="doi".

```
<front>
  <article-meta>
    <article-id pub-id-type="doi">10.2903/j.efsa.2023.7904</article-id>
```

4.2 Adoption date

The traditional foods do not have adoption date, instead they have approval date, which is not included in the JATS files and therefore cannot be extracted. The novel foods have adoption date.

front - notes -fn-group - v nejake z nich je Adopted group6 - ADOPTED: in fn-group

```
<front>
  <notes>
    <fn-group>
      <fn id="efs26305-note-1203" xml:lang="en">
        <p>Adopted: 22 October 2020</p>
      </fn>
    </fn-group>
  </notes>
</front>
```

4.3 Publication Date

Publication date can be found in article-meta - pub-date - pub-type="epub".

4.4 Question Number

question number je taky v jedne fn-group

```
<front>
  <notes>
    <fn-group>
      <fn id="efs26305-note-1002" xml:lang="en">
        <p>
          <b>Question number:</b>
          EFSA-Q-2020-00491
        </p>
      </fn>
    </fn-group>
  </notes>
</front>
```

4.5 Scientific Panels

Usually NDA, for traditional foods EFSA, sometimes also GMO.

group 3 - v title: Panel on dietetic products, nutrition and allergies group 6 - article-meta : contrib-group : collab(collab-type="authors") : NDA

look into collab - some panel there? (gmo, nda) -i added - no panel there? - look into title -i some panel there -i added -i no panel there -i panel='EFSA'

ve formátu **Group 2** je to až někde úplně dole: v body → pak najít sekci s tímto titlem

```
<sec id="efs28416-sec-0024" xml:lang="en">
  <title>QUESTION NUMBER</title>
  <p xml:lang="en">EFSA-Q-2020-00491</p>
</sec>
```

pro nalezení pouzity regex:

```
r 'EFSA[-]Q[-]\S+'
```

4.6 Scientific Officer

EFSA provided us with list of possible scientific officers. vic popsat jak toto funguje

4.7 Mandate type

-order matters

- guidance - always has guidance in the title - traditional food - always have notification of XX as traditional food
- extension of use - new dossier - the remaining ones
- nutrient source - can be added

4.8 Novel Food Name

- v title v article meta → article-title

4.9 Category

One NF can fall under multiple categories. Categories as per Regulation Article 3 of 2015/2283. Traditional food guidelines specify following categories

Categories are defined in specific european regulations. Each regulation defines its own set of categories.

REGULATION (EC) No 258 /97 from 27 January 1991

- foods and food ingredients containing or consisting of genetically modified organisms within the meaning of Directive 90 /220 /EEC;

- (b) foods and food ingredients produced from, but not containing, genetically modified organisms;
- (c) foods and food ingredients with a new or intentionally modified primary molecular structure;
- (d) foods and food ingredients consisting of or isolated from micro-organisms, fungi or algae;
- (e) foods and food ingredients consisting of or isolated from plants and food ingredients isolated from animals, except for foods and food ingredients obtained by traditional propagating or breeding practices and having a history of safe food use;
- (f) foods and food ingredients to which has been applied a production process not currently used, where that process gives rise to significant changes in the composition or structure of the foods or food ingredients which affect their nutritional value, metabolism or level of undesirable substances .

commission recommendation 97/618/EC from 29 July 1997

- (a) pure chemicals or simple mixtures from non-GM sources;
- (b) complex NF from non-GM source;
- (c) GM plants and their products;
- (d) GM animals and their products;
- (e) GM microorganisms and their products;
- (f) foods produced using a novel process.

REGULATION (EU) 2015/2283 from 25 November 2015, Article 3

- (i) food with a new or intentionally modified molecular structure, where that structure was not used as, or in, a food within the Union before 15 May 1997;
- (ii) food consisting of, isolated from or produced from microorganisms, fungi or algae;
- (iii) food consisting of, isolated from or produced from material of mineral origin;
- (iv) food consisting of, isolated from or produced from plants or their parts, except when the food has a history of safe food use within the Union and is consisting of, isolated from or produced from a plant or a variety of the same species obtained by:

- traditional propagating practices which have been used for food production within the Union before 15 May 1997; or
- non-traditional propagating practices which have not been used for food production within the Union before 15 May 1997, where those practices do not give rise to significant changes in the composition or structure of the food affecting its nutritional value, metabolism or level of undesirable substances;
- (v) food consisting of, isolated from or produced from animals or their parts, except for animals obtained by traditional breeding practices which have been used for food production within the Union before 15 May 1997 and the food from those animals has a history of safe food use within the Union;
- (vi) food consisting of, isolated from or produced from cell culture or tissue culture derived from animals, plants, micro-organisms, fungi or algae;
- (vii) food resulting from a production process not used for food production within the Union before 15 May 1997, which gives rise to significant changes in the composition or structure of a food, affecting its nutritional value, metabolism or level of undesirable substances;
- (viii) food consisting of engineered nanomaterials as defined in point (f) of this paragraph;
- (ix) vitamins, minerals and other substances used in accordance with Directive 2002/46/EC, Regulation (EC) No 1925/2006 or Regulation (EU) No 609/2013, where:
 - a production process not used for food production within the Union before 15 May 1997 has been applied as referred to in point (a) (vii) of this paragraph; or
 - they contain or consist of engineered nanomaterials as defined in point (f) of this paragraph;
- (x) food used exclusively in food supplements within the Union before 15 May 1997, where it is intended to be used in foods other than food supplements as defined in point (a) of Article 2 of Directive 2002/46/EC;

5 Composition extraction

why was deemed important how was performed - gpt assistant + query success rate

6 Implementation

-mozna neni potreba zas tak popisovat Two main classes:

JATS Opinion

- encompass the formal parts of the opinion, abstract the inner structure of the opinions into common interface
- it will abstract the formal things about the article: DOI, Title, EFSA Question, adoption date, publication date, URL, authors(panels)
- also it will abstract the different headings from different types of opinion and match them to the same structure

Opinion

- the scientific opinion itself, it will accept JATS opinion in constructor
- it will have subclass for each category, which will contain the category specific information
- it will extract NF Name, Applicant, Country of origin, NF code, type of mandate, regulation, outcome, common name, trade names, food form
- also proposed uses + target population
- availability of ADME studies
- allergenicity
- food category – for FOODS
- for **Microorganisms**: type, genus, species, strain, QPS(?)
- for **Plants**: type, common name, botanical name, genus, species, part used
- for **Animals**: genus, species, subspecies, part used
- for **Cell or Tissue**: genus, species, cell type
- for **Chemicals**: common name, IUPAC name, CAS notation, SMILES, molecular formula, InChi