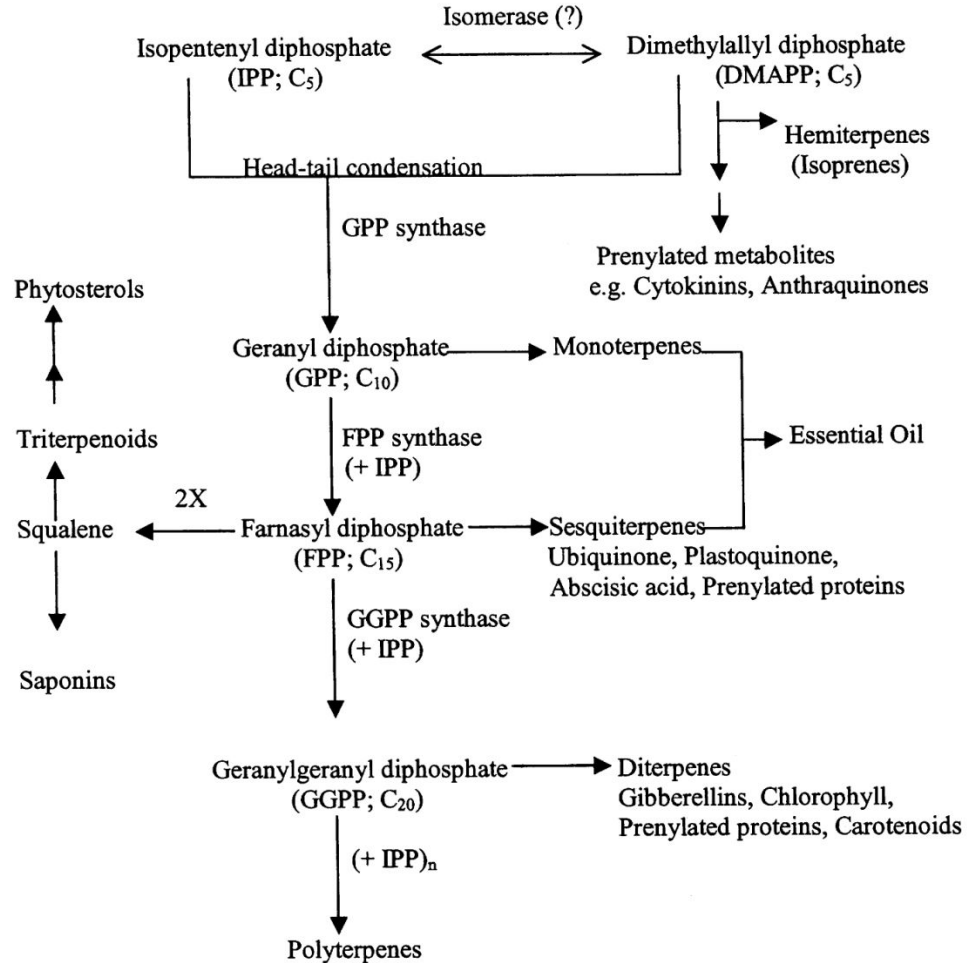


Biosynthetic pathway extraction from images

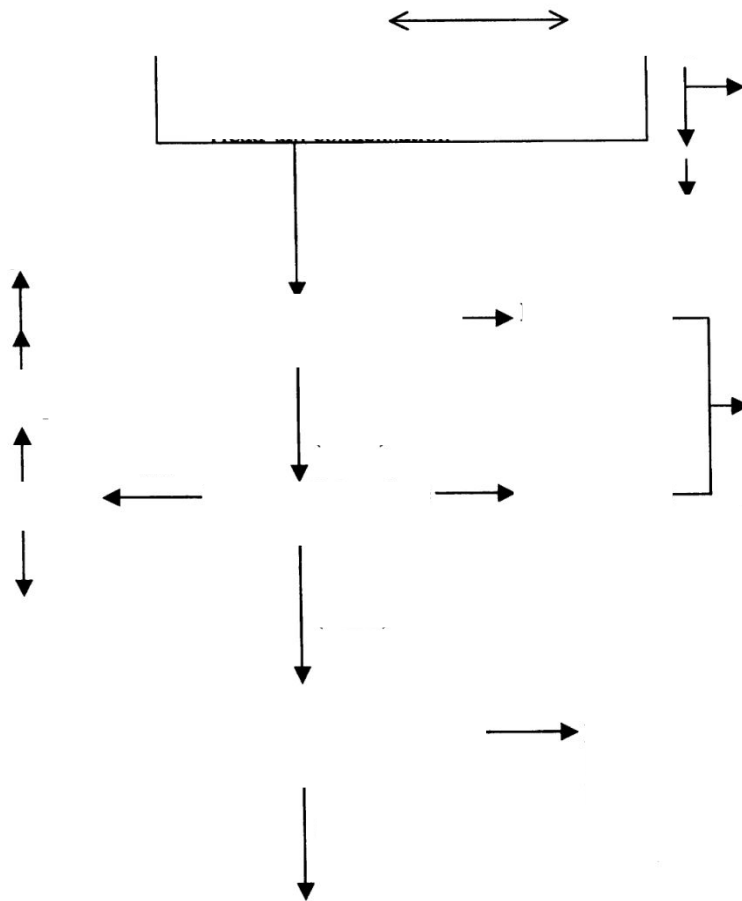
# Plan

The problem can broadly be divided into two parts:

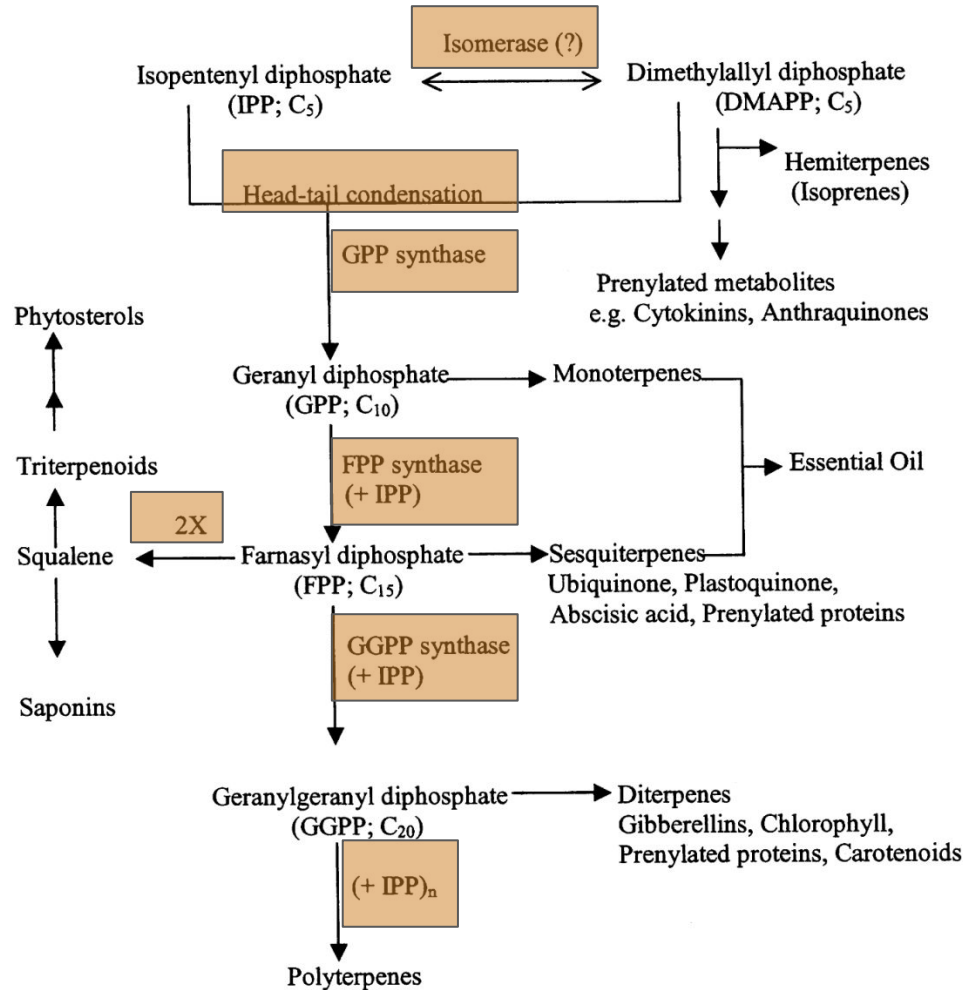
- Logical structure of the diagram
- The contextual meaning of the text



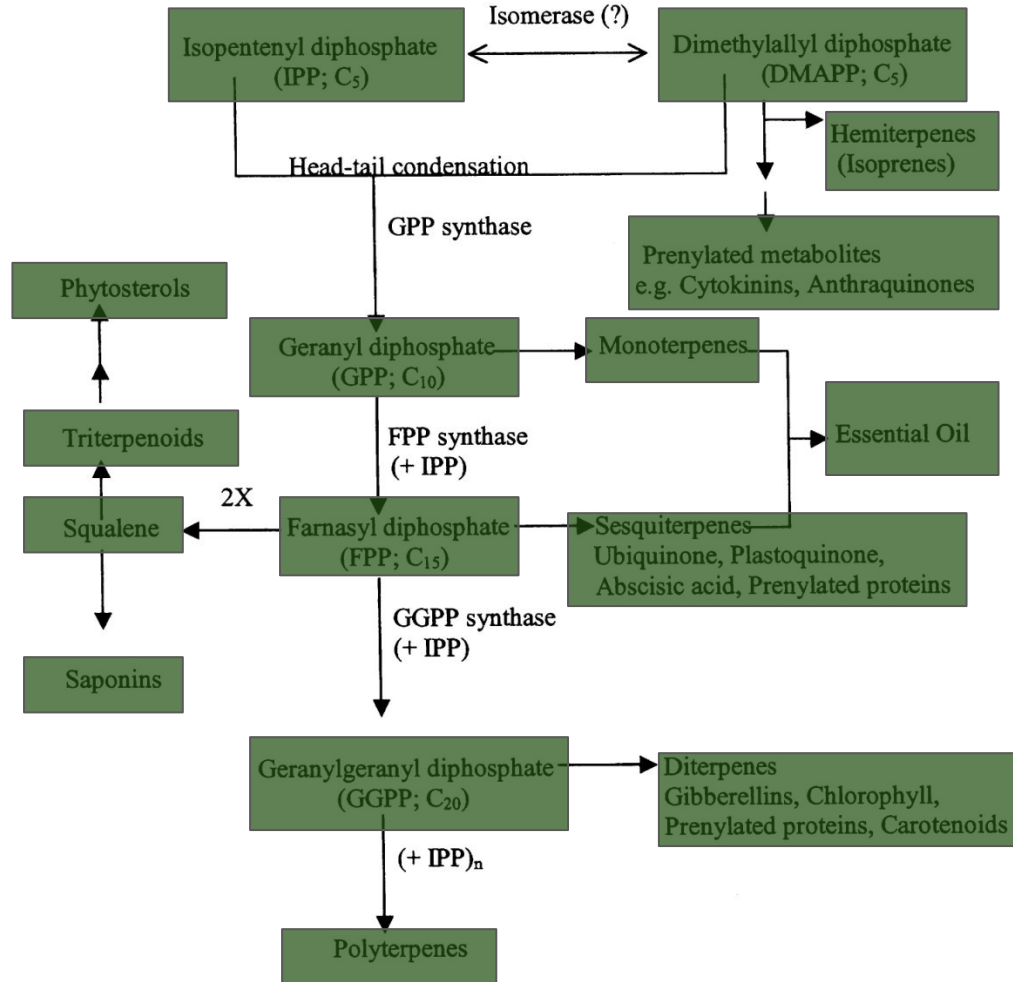
## Arrow Layer



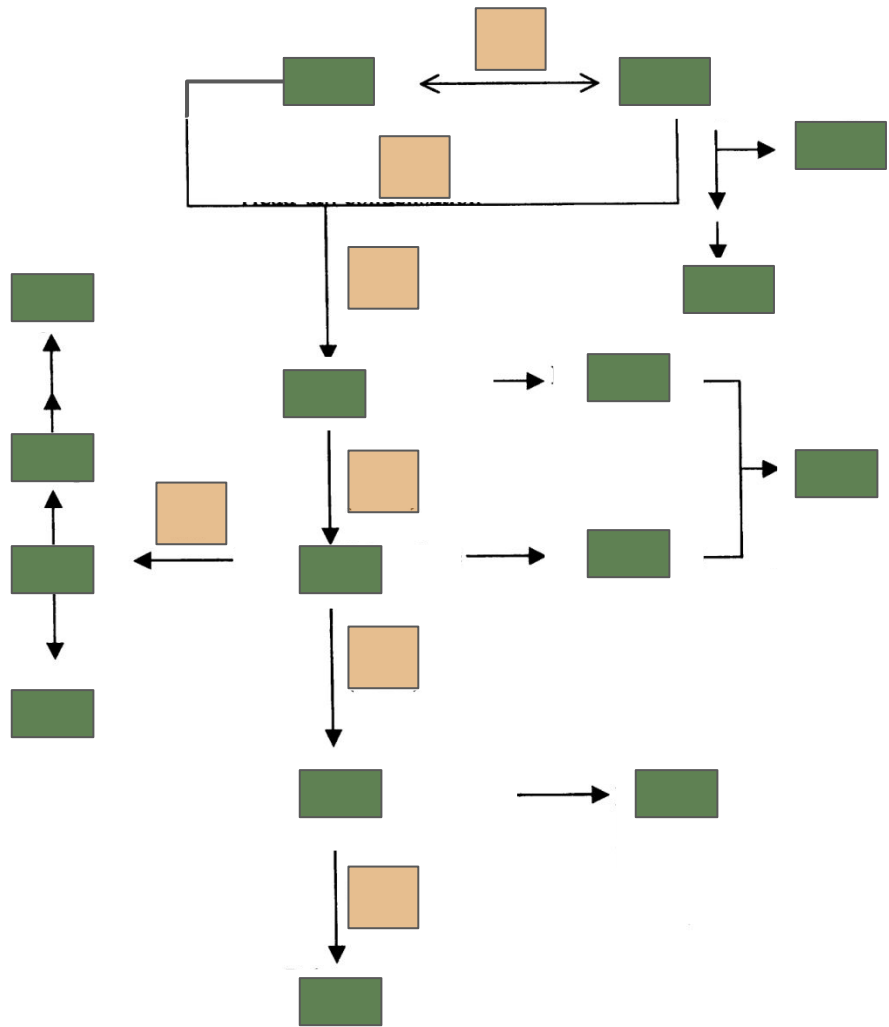
Image



## Products and Reactants Layer



Logical  
Semantics



# Challenges

- Overlapping characters and arrows
- Curved arrows
- No standard format for biosynthetic pathway diagrams



Image segment

Text output

Geranyl diphosphate—  
(GPP; C<sub>10</sub>)

Tesseract

Geranyl diphosphate (GPP; C<sub>10</sub>)

GPP synthase

Tesseract

GPP synthase

Prenylated metabolites  
e.g. Cytokinins, Anthraquinones

Tesseract

Prenylated metabolites  
e.g. Cytokinins, Anthraquinones

Lexical  
Semantics

Geranyl diphosphate—  
(GPP; C<sub>10</sub>)

Partial Chemical  
formula

Chemical  
Name

Geranyl diphosphate (GPP; C<sub>10</sub>)

Synonym/abbreviation

---

Prenylated metabolites  
e.g. Cytokinins, Anthraquinones

Chemical  
Type

Prenylated metabolites  
e.g. Cytokinins, Anthraquinones

Examples

# Challenges

- Low resolution JPEGs
- Overlapping characters
- Different kinds of information present

# Plan

The problem can broadly be divided into two parts:

- Logical structure of the diagram
- The contextual meaning of the text

Week 3 - 4:

- Resolving arrows
- Segmenting the image
- OCR of Islands of text

Week 1 - 2 (Done):

- Project planning
- Selecting target images
- Exploring image processing tools

Week 5 - 6:

- Creating diagram structure map from image segments and arrow orientations
- Annotating terms using dictionary