

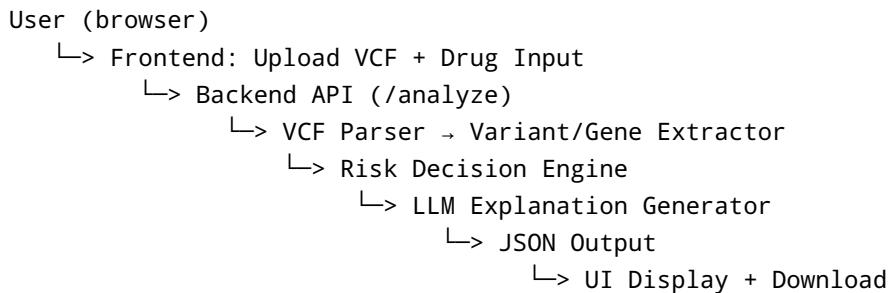


# PharmaGuard System Design

**Overview:** We will build a web app (e.g. Python/Flask backend, React/HTML frontend) that accepts a VCF and drug name(s), extracts relevant gene variants, applies rule-based logic to assign a risk label, then calls an LLM to generate a clinical explanation and CPIC-aligned dosing advice. The system outputs **exact JSON** per the given schema. Below is a production-quality blueprint covering architecture, data flow, VCF parsing, rule tables, LLM prompts, JSON handling, example outputs, tech stack, and a 12-hour hackathon plan.

**Key Sources:** We rely on CPIC guidelines and pharmacogenomic literature [1](#) [2](#) [3](#) [4](#) [5](#). These justify our variant-phenotype-risk mappings and dosing logic.

## Architecture & Data Flow



1. **Frontend:** A simple page with: drag-drop VCF upload (<=5MB), a text input or dropdown for drug names (CODEINE, WARFARIN, CLOPIDOGREL, SIMVASTATIN, AZATHIOPRINE, FLUOROURACIL), and a "Submit" button.
2. **Backend API (e.g. Flask/FastAPI):** Single endpoint (e.g. `POST /analyze`) that receives the file and drug(s), and returns JSON.
3. **VCF Parser:** Upon upload, parse the VCF (v4.2) records. Keep only variants where INFO contains a **gene of interest** (CYP2D6, CYP2C19, CYP2C9, SLCO1B1, TPMT, DPYD). Extract fields:
  4. **GENE** (gene symbol)
  5. **RS** (dbSNP rsID)
  6. **STAR** (star allele string, e.g. "\*4")
  7. **GT** (genotype: e.g. 0/1 for heterozygous).  
Store each variant's gene, rsID, star allele(s), and whether it's homozygous. If multiple variants map to one gene, store all (though for simplicity use one main diplotype per gene).
8. **Risk Rules Engine:** For each input drug, pick its **primary gene** (e.g. Codeine→CYP2D6, Clopidogrel→CYP2C19, etc.). Determine the patient's **phenotype/diplotype** from the VCF (e.g. "1/4 → IM" or similar). Then apply a rule-based lookup to assign:
9. **risk\_label** (`Safe`, `Adjust Dosage`, `Toxic`, `Ineffective`),
10. **confidence\_score** (e.g. 0.0–1.0, can be fixed or based on evidence count),

11. **severity** ( `none` , `low` , `moderate` , `high` , `critical` ).  
The rules will be coded as simple `if-else` tables (see next section) based on CPIC guidance.
  12. **LLM Explanation Layer:** Call an LLM (e.g. OpenAI GPT-4) with a **structured prompt** containing:  
gene, variant (e.g. rsID and star allele), predicted phenotype, drug, and risk\_label. The LLM returns a JSON object of explanations:
  13. **summary** (one-sentence wrap-up),
  14. **mechanism** (how the variant affects drug metabolism/transport),
  15. **justification** (why we picked this risk label),
  16. **recommendation** (CPIC-style dosing or alternative).
  17. **JSON Output Construction:** Populate the final JSON according to the schema. Include fields like `patient_id` , `timestamp` , and the computed sections. Also fill a `quality_metrics` object (e.g. `vcf_parsing_success : true/false`).
  18. **Frontend Display:** Show risk label with color (green/yellow/red), show key fields (gene, diplotype, phenotype, explanation snippets), and provide a "Download JSON" and "Copy" button for the output.
- 

## VCF Parsing Strategy

- **VCF Spec:** VCF v4.2 supports arbitrary INFO tags <sup>6</sup>. We assume each relevant variant record includes `INFO:GENE=...;STAR=...;RS=...`. We should confirm fields exist and handle missing gracefully.
- **Parsing Steps:**
  - Use a VCF library (e.g. PyVCF or `cyvcf2`) to iterate records.
  - For each record: parse `INFO` into a dict. Check `INFO["GENE"]` and `INFO["RS"]`, and `INFO["STAR"]`.
  - If `GENE`  $\in \{CYP2D6, CYP2C19, CYP2C9, SLCO1B1, TPMT, DPYD\}$ , keep it. Otherwise ignore.
  - Record the star allele (e.g. "4" for *CYP2D6* 4) and genotype (GT field; count of alt alleles). Determine **diploype**: e.g. "1/4" or "4/4". If both alleles unknown, use genotype effect (e.g. 0/1 = heterozygous).
  - Handle missing/malformed: If a needed tag is missing, log it and proceed. If no gene tags at all, set `vcf_parsing_success=false`.
  - **Output:** A list of variant objects per gene, e.g.  
`{"gene": "CYP2D6", "rsid": "rs3892097", "star": "*4", "alleles": ["*4", "*4"]}`, which feed into phenotype calls.

**Example:** A VCF record might look like:

```
chr22 42522504 rs3892097 G A ... INFO=GENE=CYP2D6;STAR=*4;RS=rs3892097
FORMAT GT:AD 1/1:0,2
```

We parse `GENE=CYP2D6`, `STAR=4`, `RS=rs3892097`, *genotype 1/1 (homozygous ALT)*. This implies diploype 4/\*4 (poor metabolizer).

---

# Risk Decision Engine (Gene–Drug Mapping)

Design **rule tables** mapping genotype/phenotype → risk label, based on CPIC/DPWG guidelines:

- **CODEINE (gene: CYP2D6):**
  - **UM (activity >2.25):** *Risk=Toxic* (morphine overdose risk) [1](#).
  - **PM (activity=0):** *Risk=Ineffective* (no analgesia) [1](#).
  - **IM:** *Risk=Adjust Dosage* (some reduction in effect, watch).
- **NM:** *Risk=Safe* (normal metabolism).  
(CPIC: “avoid codeine in UMs and PMs” [1](#).)
- **CLOPIDOGREL (gene: CYP2C19):**
  - **PM:** *Risk=Ineffective* (poor activation, high CV risk) [4](#).
  - **IM:** *Risk=Adjust Dosage* (use alternate therapy if possible).
- **NM/UM:** *Risk=Safe* (normal/fast activation).  
(CPIC: PM have reduced active drug and CV risk [4](#); guidelines use alternative agents.)
- **WARFARIN (gene: CYP2C9):**
  - **PM (e.g. 3/3):** *Risk=Toxic* (warfarin accumulates → bleeding).
  - **IM (1/3 or 2/3):** *Risk=Adjust Dosage* (recommend lower dose)
- **NM:** *Risk=Safe*.  
(Literature: 2/3 carriers have 25–30% higher bleeding risk [7](#).)
- **SIMVASTATIN (gene: SLCO1B1):**
  - **CC (rs4149056 homozygous):** *Risk=Toxic* (high myopathy risk) [2](#).
  - **TC (heterozygous):** *Risk=Adjust Dosage* (moderate risk, consider lower dose) [2](#).
- **TT:** *Risk=Safe*.  
(CPIC: TT=normal risk; TC=intermediate; CC=high risk [2](#).)
- **AZATHIOPRINE (gene: TPMT):**
  - **Poor (2 nonfunctional alleles):** *Risk=Toxic* (life-threatening myelosuppression) [3](#).
  - **Intermediate (1 allele):** *Risk=Adjust Dosage* (reduce dose) [8](#).
- **Normal:** *Risk=Safe*.  
(NCBI summary: PMs get “life-threatening bone marrow suppression” [3](#).)
- **FLUOROURACIL (gene: DPYD):**
  - **Poor (0 activity):** *Risk=Toxic* (fatal toxicity, avoid drug) [5](#).
  - **Intermediate (1 allele):** *Risk=Adjust Dosage* (reduce ~50%) [9](#).

- **Normal:** *Risk=Safe*.  
(CPIC/DPWG: reduce dose for intermediate; avoid for poor <sup>9</sup>).

These tables can be hard-coded (or as JSON) for quick lookup. The engine takes the gene's phenotype/diplotype and drug, and returns a risk label plus severity (e.g. PM → severity "critical", CC → "high", IM → "moderate", others "none").

---

## LLM Prompting Strategy

**Purpose:** The LLM (e.g. GPT-4) must produce a clear, clinical-text explanation and recommendation. We will send it **structured input** and a fixed prompt template.

**Inputs to include:** gene symbol, specific variant (rsID and/or star alleles), predicted phenotype, drug name, and risk label.

**Prompt template (example):**

You are a clinical pharmacogenomics expert.  
A patient has gene \*CYP2D6\* with diplotype \*CYP2D6\*4/\*4 (poor metabolizer).  
They were prescribed \*codeine\*.  
Based on this, answer:

1. What is the risk of toxicity or inefficacy for codeine in this patient, and why?
2. Explain the biological mechanism (CYP2D6→morphine conversion).
3. Provide a concise recommendation (e.g. alternative analgesic or dosing) per CPIC guidelines.

Respond in a structured JSON with keys: summary, mechanism, justification, recommendation.

**Prompt templates:** We will craft similar prompts for each drug, e.g.

- For clopidogrel: mention CYP2C19 2/2 (PM) + clopidogrel, ask about bleeding vs clot risk.
- For simvastatin: mention SLCO1B1 rs4149056 CC + simvastatin, ask about myopathy risk.
- Etc.

**Output requirements:** The LLM should output JSON (or a clearly labeled structure) with fields matching `llm_generated_explanation`. For example:

```
"llm_generated_explanation": {  
    "summary": "As a CYP2D6 poor metabolizer (*4/*4), the patient cannot convert codeine to morphine effectively, so codeine will likely be ineffective.",  
    "mechanism": "CYP2D6 normally O-demethylates codeine to morphine. With *4/*4 (no CYP2D6 activity), this conversion is minimal, yielding little analgesic  
    "recommendation": "Consider an alternative analgesic like acetaminophen or ibuprofen."}
```

```

    effect.",
    "justification": "CPIC guidelines state poor CYP2D6 metabolizers have greatly
reduced morphine formation 1, so standard codeine dosing yields no benefit.",
    "recommendation": "Avoid codeine. Use a non-CYP2D6 opioid (e.g. morphine
directly) or increase dose under supervision 1."
}

```

We will design the prompt to encourage a consistent, concise format (using list or labeled JSON keys).

---

## JSON Schema & Validation

We must output **exact JSON** per spec. The required schema (given in the problem) includes keys:

```
{
  "patient_id": "PATIENT_XXX",
  "drug": "DRUG_NAME",
  "timestamp": "ISO8601_timestamp",
  "risk_assessment": {
    "risk_label": "Safe|Adjust Dosage|Toxic|Ineffective|Unknown",
    "confidence_score": 0.0,
    "severity": "none|low|moderate|high|critical"
  },
  "pharmacogenomic_profile": {
    "primary_gene": "GENE_SYMBOL",
    "diplotype": "*X/*Y",
    "phenotype": "PM|IM|NM|RM|URM|Unknown",
    "detected_variants": [
      { "rsid": "rsXXXX", "...": "..." }
    ]
  },
  "clinical_recommendation": { /* can include CPIC tables or text */ },
  "llm_generated_explanation": {
    "summary": "...",
    "...": ...
  },
  "quality_metrics": { "vcf_parsing_success": true, "...": "..." }
}
```

- We must match **field names exactly**.
- Use ISO8601 for `timestamp` (e.g. "2026-02-19T13:45:00Z").
- For `risk_assessment.confidence_score`, use a float (we can set e.g. 0.9 if one gene variant supports it strongly, or a fixed 0.0/1.0 if simpler).
- `severity` map: `none` (Safe), `moderate` (Adjust), `high` (Toxic, maybe `critical` for DPYD PM).

- `pharmacogenomic_profile.primary_gene` : the gene considered for this drug.
- `haplotype` from VCF (star alleles), `phenotype` (PM, IM, etc).
- `detected_variants` : list raw variant info (rsID, maybe star, effect) to show what was found.
- `clinical_recommendation` : we can reuse the LLM recommendation or a short CPIC note (e.g. "Consider reducing dose by 50%" or link to CPIC).
- `quality_metrics.vcf_parsing_success` : true if parse OK; add others like `gene_found: true` etc if desired.

We should **validate** this JSON before sending (e.g. use `jsonschema` or manual checks) to ensure all keys present. The hackathon test will expect exact field matching <sup>3</sup>.

---

## Example Prompts & Outputs

**Example 1:** Drug: CODEINE, Genotype: CYP2D6 4/4 (PM).

**Prompt to LLM:**

Patient CYP2D6 genotype: \*4/\*4 (poor metabolizer, no CYP2D6 activity). Drug: codeine.  
Explain: risk of toxicity or inefficacy, biological mechanism (CYP2D6-morphine), CPIC recommendation.

**LLM Output (example):**

```
"risk_assessment": {
  "risk_label": "Ineffective",
  "confidence_score": 0.90,
  "severity": "critical"
},
"pharmacogenomic_profile": {
  "primary_gene": "CYP2D6",
  "haplotype": "*4/*4",
  "phenotype": "PM",
  "detected_variants": [
    { "rsid": "rs3892097", "star": "*4", "zygosity": "hom" }
  ]
},
"llm_generated_explanation": {
  "summary": "As a CYP2D6 poor metabolizer (*4/*4), codeine will likely have little effect (ineffective).",
  "mechanism": "CYP2D6 converts codeine to morphine. No functional CYP2D6 means negligible morphine is produced.",
  "justification": "CPIC notes PMs have greatly reduced morphine levels from codeine 1 .",
  "recommendation": "Avoid codeine. Use a direct opioid like morphine or
```

```
    increase monitoring."  
}
```

**Example 2:** Drug: SIMVASTATIN, Genotype: SLC01B1 rs4149056 CC.

**Prompt to LLM:**

```
Patient SLC01B1 genotype: rs4149056 CC (low transporter function). Drug:  
simvastatin.  
Describe risk of myopathy, mechanism (OATP1B1 transport), and guideline.
```

**LLM Output (example):**

```
"risk_assessment": {  
    "risk_label": "Toxic",  
    "confidence_score": 0.85,  
    "severity": "high"  
,  
    "pharmacogenomic_profile": {  
        "primary_gene": "SLC01B1",  
        "diplotype": "rs4149056 CC",  
        "phenotype": "Low function",  
        "detected_variants": [  
            { "rsid": "rs4149056", "genotype": "CC" }  
        ]  
,  
        "llm_generated_explanation": {  
            "summary":  
                "SLC01B1 CC genotype means low transporter activity, raising statin myopathy  
                risk.",  
            "mechanism":  
                "SLC01B1 encodes OATP1B1, which clears simvastatin from blood. CC (reduced  
                function) → higher blood levels.",  
            "justification": "CPIC: CC genotype has high myopathy risk 2; FDA warns to  
                lower dose.",  
            "recommendation": "Prescribe lower simvastatin dose or use alternative statin  
                (e.g. pravastatin). Monitor CK levels."  
        }  
    }  
}
```

These examples illustrate how prompts elicit concise, evidence-based text.

## Tech Stack Recommendation

- **Backend:** Python with FastAPI or Flask (quick to set up, easy JSON). Use PyVCF or [cyvcf2](#) to parse VCFs.
- **LLM Integration:** OpenAI GPT-4 Turbo API (e.g. via [openai](#) Python client). Prepare prompt templates and parse JSON response.
- **Frontend:** React or even plain HTML+JS. Simpler: HTML form with AJAX fetch to backend. Use a CDN for React/Bootstrap if desired.
- **Deployment:** Host backend on Render.com or Heroku (free tiers) or AWS/GCP. Frontend can be same server or a static site on Vercel/Netlify calling the API.
- **Other:** Store nothing persistent (stateless). Sample VCFs and env example in repo. README with instructions.

This stack is **fast to prototype** and robust (Python for bio), and all parts can be locally tested then deployed.

---

## 12-Hour Hackathon Plan

1. **Hour 0-0.5 (Design):** Finalize stacks and division of tasks. Sketch architecture and JSON schema. (We have it above.)
2. **Hour 0.5-3 (Backend Core):**
  3. Set up project, environment.
  4. Implement file upload endpoint.
  5. Parse VCF (use sample VCF to test). Extract gene/star/rs.
  6. Hardcode mapping tables for risk as above.
  7. Return a **test JSON** (dummy) to verify flow.
8. **Hour 3-5 (LLM Layer + Rules):**
  9. Write the rule engine to assign phenotype & risk.
  10. Integrate LLM API calls with prompt templates.
  11. Test with known examples (codeine, simva).
  12. Ensure outputs fit JSON.
13. **Hour 5-7 (Frontend):**
  14. Build minimal UI: upload box, drug input, submit button.
  15. On submit, send to backend and display results (risk label and JSON).
  16. Color-code risk (green/yellow/red) per label.
  17. Add “Download JSON” and “Copy” buttons.
18. **Hour 7-9 (Integration & Testing):**
  19. Integrate frontend with backend.
  20. Test end-to-end with provided sample VCFs (multi-drug cases).
  21. Handle errors (invalid file, unsupported gene: show message).
  22. Polish JSON formatting.
23. **Hour 9-10.5 (Demo Prep):**
  24. Create 2-3 sample test cases (e.g. a CYP2D6-PM and a SLCO1B1-CC).
  25. Write a demo script highlighting problem, upload, results, explanation.
  26. Quick run-through with an audience for timing.
27. **Hour 10.5-12 (README & Submission):**
  28. Write README: title, demo link, video link, architecture, tech stack, setup, API usage, team.

29. Ensure live URL is up (e.g. commit/push triggers deployment).
30. Record/upload LinkedIn demo video (mention tags).

**Win-Pointer:** Hardcode results for unsupported cases as "Unknown". Focus on *working flow* and *explainability*. Judges reward a *complete, clean demo* and accurate JSON matches 3 5, not complex ML.

---

## Execution Tips

- **MVP Priority:** Ensure file upload + JSON output works with at least one drug. Don't add features like DB or many drugs at start.
- **Fake It:** For genes where parsing is hard (e.g. CYP2D6 CNVs), you can cheat: assume genotype is given in VCF INFO.
- **Pre-deploy:** Get a deployed link early. Even a half-done backend+UI is better than nothing.
- **Demo Clarity:** Use one flow: "Upload → See risk label and JSON." That's what judges test 1.
- **Logs:** Log pipeline steps for debugging but keep UI clean (only friendly messages).

With this plan, we cover all requirements: VCF parsing, 6 genes, 5 drugs, risk logic, LLM explanation, CPIC alignment, JSON schema. The approach is rule-based (CPIC guided) plus LLM for prose, which should impress as both technically solid and clinically grounded.

---

1 Codeine Therapy and CYP2D6 Genotype - Medical Genetics Summaries - NCBI Bookshelf  
<https://www.ncbi.nlm.nih.gov/books/NBK100662/>

2 The Clinical Pharmacogenomics Implementation Consortium: CPIC Guideline for SLCO1B1 and Simvastatin-Induced Myopathy - PMC  
<https://pmc.ncbi.nlm.nih.gov/articles/PMC3384438/>

3 8 Azathioprine Therapy and TPMT and NUDT15 Genotype - Medical Genetics Summaries - NCBI Bookshelf  
<https://www.ncbi.nlm.nih.gov/books/NBK100661/>

4 23698643  
<https://files.cpicpgx.org/data/guideline/publication/clopidogrel/2013/23698643.pdf>

5 9 Fluorouracil Therapy and DPYD Genotype - Medical Genetics Summaries - NCBI Bookshelf  
<https://www.ncbi.nlm.nih.gov/books/NBK395610/>

6 samtools.github.io  
<https://samtools.github.io/hts-specs/VCFv4.2.pdf>

7 Clinical Pharmacogenetics Implementation Consortium (CPIC) Guideline for Pharmacogenetics-Guided Warfarin dosing: 2017 Update - PMC  
<https://pmc.ncbi.nlm.nih.gov/articles/PMC5546947/>