



Improving Clinical Utility of ICD and DSM with an AI-Powered Symptom Diagnosis Tool

Serious Problems with ICD/DSM Utility: Both the **International Classification of Diseases (ICD)** and the **Diagnostic and Statistical Manual of Mental Disorders (DSM)** have long been criticized for poor **clinical utility**, meaning they often fail to be truly useful in day-to-day patient care ¹. As psychologist Geoffrey Reed noted, *“Serious problems with the clinical utility of both the ICD and the DSM are widely acknowledged.”* ². In practice, clinicians and researchers find that these rigid categorical systems **don’t always match real-world patient presentations or aid treatment decisions**. For example, a 2023 analysis stated that the DSM/ICD categories “do not fit either the clinical presentations of patients, recent discoveries from genetics and neurobiology, nor therapeutic choices well” ³. In other words, patients’ symptoms often span multiple categories or don’t neatly fit the definitions, making the classifications less helpful at the bedside. There is also a **lack of alignment with how clinicians naturally think** about illnesses: studies show that experienced clinicians tend to mentally group disorders differently than the official taxonomy does ⁴. The result is that many professionals find the ICD/DSM cumbersome – with **too many diagnostic codes, overlapping labels, and insufficient guidance for treatment** ⁵. All these issues undermine the utility of these systems in clinical practice.

Impact on Care: The stakes for improving diagnostic tools are high. Clinical utility isn’t a trivial matter – it *“affects the daily lives of practitioners and is also a global public health issue”* ². When classification systems are unwieldy or imprecise, **misdiagnosis and treatment delays** become more likely. In mental health, Reed notes that most people with mental disorders worldwide receive no treatment, partly due to under-identification ⁶. In general medicine, diagnostic errors are distressingly common; roughly *10-15% of medical diagnoses may be wrong* ⁷. Such errors carry a huge human and economic cost – by one estimate, misdiagnoses cost the U.S. healthcare system up to **\$100 billion per year** ⁸. Clearly, better tools for identifying the correct diagnosis (and doing so early) could save lives and resources. A more **clinically useful system** would help practitioners recognize conditions more accurately and communicate clearly, ultimately improving patient outcomes ⁶. This is why the **WHO’s revision for ICD-11** explicitly prioritized clinical utility, defining it in terms of (a) effective communication of information, (b) ease of use and “goodness of fit” to patient presentations, and (c) usefulness in guiding treatment decisions ⁹. Any innovation that boosts these aspects – that makes diagnoses easier, more accurate, and more actionable – would be a game-changer for healthcare.

Proposed Solution – AI Symptom-to-Diagnosis MCP Server: To tackle these issues, we propose a **medical symptom diagnosis tool built as an MCP server** (Model Context Protocol server) for the Hugging Face *Agents & MCP Hackathon 2025 (Track 1)*. In essence, this is an **AI-powered backend service (built with Gradio)** that takes in **patient-entered symptoms in natural language** and outputs a structured list of probable diagnoses **mapped to ICD-10 codes**, complete with confidence scores. The system leverages a **local medical knowledge base** and advanced language model reasoning to bridge the gap between the patient’s story and the formal diagnostic categories. The choice of ICD-10 coding ensures outputs are in a globally recognized format (facilitating communication and billing), while the design is extensible to ICD-11 or DSM criteria in the future. By packaging this in an MCP server, any AI agent or application can query it –

for example, a chatbot could call the MCP server with a patient's symptom description and receive a JSON response of likely diagnoses.

- **Natural Language Input – Ease of Use:** A key advantage of this tool is **usability**. Instead of forcing clinicians or patients to navigate giant codebooks or memorize DSM criteria, the **user can simply describe symptoms in plain language** (e.g. *"persistent sad mood, loss of interest, and insomnia"*). The AI will interpret and map this description to relevant diagnostic entities. This addresses the **"ease of use" component of clinical utility** – making the process far more intuitive. It aligns with how people actually talk about their problems, thus *"representing the way that illness is experienced by the person"* (an often missing element in DSM/ICD) ¹⁰. By catering to patient language, the tool lowers the barrier for diagnosis and ensures important details aren't lost in translation. For hackathon judges, this showcases a **user-friendly interface**: a simple input box for symptoms and a clear output, demonstrating Gradio's capability to handle natural text inputs and produce useful results.
- **Knowledge Base Mapping – Goodness of Fit:** Under the hood, the MCP server uses a **curated knowledge base of symptom-disease relationships** (drawn from clinical guidelines and databases) to ensure that the suggested diagnoses are grounded in medical evidence. This helps address the **"goodness of fit" problem** with current classifications. In practice, clinicians often complain that diagnoses feel like forcing a square peg into a round hole – patients rarely match textbook definitions exactly. Our AI approach instead **considers multiple possible diagnoses** and rates how well each fits the reported symptoms. For example, a patient's symptoms *"fatigue, weight loss, thirst"* might fit **diabetes** strongly but also hint at **hyperthyroidism** or **chronic anemia**. The MCP server would return all these possibilities with confidence scores (e.g. Diabetes Mellitus – ICD-10 E11 – 90% likely; Hyperthyroidism – E05 – 40%; Iron-deficiency anemia – D50 – 30%). By presenting a **differential diagnosis** list, the tool acknowledges uncertainty and comorbidity, reflecting reality better than a single rigid label. This directly mitigates the issues of the DSM/ICD categorical approach by introducing **probabilistic, nuanced assessments**. Notably, providing **confidence scores** is a practical innovation – clinicians can gauge which diagnoses are most likely yet remain aware of alternatives. This feature impressed our clinical advisors as something standard manuals don't offer, but which is very valuable in decision-making.
- **Structured Output (JSON) – Communication & Integration:** The output is formatted as a **JSON object** containing each probable diagnosis with its ICD code and a confidence metric. For example:

```
{
  "input_symptoms": "persistent sad mood, loss of interest, insomnia",
  "likely_diagnoses": [
    {"icd10_code": "F33.1", "diagnosis": "Major Depressive Disorder, moderate",
    "confidence": 0.85},
    {"icd10_code": "F41.1", "diagnosis": "Generalized Anxiety Disorder",
    "confidence": 0.30}
  ]
}
```

This structured approach greatly improves **communication**. It means the results are machine-readable and can be easily integrated into electronic health records, decision support systems, or other agents. In terms

of clinical utility, this tackles the need for **effective information sharing** ⁹. A diagnosis in JSON with a standard code can be universally understood by other software and professionals, avoiding the ambiguity of free-text notes. During the hackathon demo, we will show how an agent (or even a simple web app) can call our MCP server API and get back this JSON, then perhaps automatically provide the patient with next steps (e.g. "It looks like you may have depression; consider seeking a mental health evaluation."). By adhering to the **Model Context Protocol**, our tool becomes a modular component – any future AI **agent** could use it as a plug-and-play medical brain. This highlights to judges that we are leveraging **Hugging Face's cutting-edge MCP architecture** properly, creating a tool that's not just a one-off demo but part of a larger **agentic ecosystem**.

Solving ICD/DSM Pain Points: Overall, this AI-powered solution directly addresses the acknowledged problems with ICD and DSM:

- **Simplifying Complexity:** Instead of navigating hundreds of categories, users get a **focused list of relevant diagnoses**. Researchers have suggested reducing the number of categories to improve clinical utility ⁵, and our tool effectively does this by filtering out irrelevant diagnoses and highlighting the few that matter for a given case. This makes the diagnostic process more **manageable and clinician-friendly**.
- **Aligning with Clinical Reasoning:** The tool mirrors how clinicians approach cases – by considering a **differential diagnosis** rather than jumping to a single label. Traditional DSM/ICD usage often pressures clinicians to pick one label per patient visit, even if they're uncertain. Our approach instead **embraces uncertainty and comorbidity**, aligning with real clinical reasoning. This can improve diagnostic accuracy and ensure less "fitting a square peg in a round hole." For example, if a patient's symptom profile doesn't perfectly match any one disorder, the AI might assign ~50% probability to two possible diagnoses – signaling the need for further evaluation or perhaps indicating the patient meets partial criteria for multiple conditions. This granularity is a step toward more **dimensional (rather than purely categorical) assessment**, which is exactly the direction experts have been advocating for improving validity and utility ³.
- **Enhancing Communication:** Because the output uses standard ICD-10 codes, it provides a **common language** for healthcare providers. One issue previously was that many U.S. clinicians train on DSM terms and may not know the equivalent ICD codes by heart ¹¹. Our tool automates the mapping from symptom description to the proper ICD code, ensuring that nothing is lost in translation. This can prevent errors in billing or epidemiological recording that occur when clinicians use unofficial diagnoses or miscode something. In the hackathon context, demonstrating **accurate ICD-10 mapping** shows attention to real-world detail (a nod to medical industry needs, which judges will appreciate).
- **Informing Treatment and Next Steps:** A diagnosis is only useful if it guides what to do next. DSM/ICD manuals themselves don't give treatment advice, and that's part of their utility problem ³. While our hackathon prototype is primarily a diagnostic tool, it could be extended to suggest **evidence-based next steps**. For instance, after listing likely diagnoses, the system could (in future iterations) pull in guideline recommendations or prompt the user to seek specific care (e.g. see an endocrinologist if diabetes is likely). Even at present, by ranking diagnoses, the tool helps a clinician prioritize which condition to evaluate or treat first. This *"usefulness for making clinical management decisions"* is a core aspect of clinical utility ⁹. We plan to mention in our presentation how such a

system could be augmented with treatment algorithms or referral suggestions, making it not just a diagnostic aid but a component in a full clinical decision support pipeline.

Why This Will Impress Hackathon Judges: Our MCP server project is innovative and ambitious, yet clearly rooted in solving a **real, widely-acknowledged problem** in healthcare. It stands out because it combines **cutting-edge AI** with a high-impact domain:

- *Technical Creativity:* We leverage the **Hugging Face MCP framework** in a novel way – turning a Gradio app into a **medical AI service**. This showcases the flexibility of MCP (if a Gradio app can serve medical diagnostics, it proves MCP's power). We also integrate a **domain-specific knowledge base** with an LLM, which is an excellent demonstration of how to ground language models in expert data for higher reliability. Judges will see that we've gone beyond a vanilla chatbot; we built a structured, domain-adapted agent tool.
- *Real-World Impact:* The project directly targets the pain points described in academic research and practitioner feedback. We can cite how **WHO and researchers have called for improved clinical utility** ², and then show our solution as a concrete step toward that goal. This alignment of a hackathon project with global health objectives and literature is likely to impress judges – it's not just a toy app, but something with purpose and potential to help millions. Healthcare is a massive sector, and AI solutions here can have outsized benefits, which adds weight to our project. (*As an aside, investors and companies are pouring resources into AI health tools – the judges will recognize that this could be more than a hackathon demo.*)
- *Innovation in User Experience:* By focusing on **patient-entered symptoms**, we emphasize accessibility. Our design considers *end-users* (patients and front-line clinicians) rather than assuming AI is only for data scientists. This user-centric approach – intuitive symptom input and clear output – follows best practices of product design in health tech. It also aligns with the Hackathon's spirit of building applications that are **practically useful and usable**. We might even do a quick live demo: e.g., input a sample symptom scenario and show the JSON result populating. Seeing an end-to-end working prototype that feels like something you could actually use gives us an edge in the competition.
- *Broad Utility and Extensibility:* Because we output standard codes and follow open protocols, the project has **future extensibility**. We can point out that this MCP server could be plugged into telehealth platforms, symptom checker websites, or used by insurance for preliminary triage. The Hugging Face community could build on it (it will be open-sourced as a Space), perhaps improving the model with more data or extending it to cover **DSM-5 diagnostic criteria for mental health** explicitly. In short, it's a foundation for a larger system. This open-ended potential is attractive in hackathons – it shows that our idea has *life beyond the one-week contest*.

Path to Winning and Beyond: By addressing a clearly identified problem (ICD/DSM utility) with an **AI-driven solution**, our project hits the key criteria: originality, technical excellence, and impact. The hackathon judges (the immediate target users for our presentation) will appreciate the robust citations and problem understanding we bring – we aren't just coding for the sake of it, we're solving something that the **WHO, APA, and researchers have flagged for years** ² ³. We will make this case in our project README and demo, showing how each feature of our MCP server maps to a pain point in clinical diagnosis. Additionally, we'll emphasize testing and accuracy; for instance, we plan to mention performance

benchmarks of symptom checkers (which on average only ~50% accuracy for top-3 diagnoses ¹²) and how our approach using a state-of-the-art model could push that higher. Even if we don't have time for a full evaluation during the hackathon, this comparative angle signals that our solution could outperform existing tools – a hint at **innovation**.

Finally, looking beyond the hackathon, this project has genuine commercial and social potential. **Winning the hackathon** would validate the concept, and the next step could be deploying it as a standalone app or API service for clinics and consumers. The demand is certainly there: millions of people search online for symptoms monthly, and clinicians worldwide struggle with diagnosis coding. A successful implementation could attract healthcare providers or startups (e.g. for licensing or partnership), aligning with our goal to **not only win but also eventually monetize** the idea. In summary, our MCP server for symptom-based diagnosis exemplifies how modern AI can “*improve the clinical utility*” of diagnostic systems ⁶ . It transforms the static, often unwieldy ICD/DSM taxonomy into a **dynamic, interactive assistant** that makes diagnosis more accurate, efficient, and user-friendly. This is a compelling narrative for the hackathon judges, and more importantly, a step toward better healthcare delivery in the real world.

Sources:

- Reed, G. M. (2010). *Toward ICD-11: Improving the clinical utility of WHO's International Classification of mental disorders*. **Professional Psychology: Research and Practice**, **41**(6), 457–464 ¹³ ⁶ . (Highlights the widely acknowledged shortcomings of ICD/DSM clinical utility and the importance of improving it.)
- Cohen, B. & Öngür, D. (2023). *Rethinking categorical psychiatric diagnoses in light of genetics and neuroscience findings*. Quoted in Aftab, A. (2023) **Psychiatry at the Margins** ³ . (Notes that current DSM/ICD categories often don't align with patient reality or biological data.)
- Flanagan, E. H., & Blashfield, R. K. (2010). *Increasing Clinical Utility by Aligning the DSM and ICD With Clinicians' Conceptualizations*. **Prof. Psychology: Research and Practice**, **41**(6), 474–481 ⁴ ⁵ . (Suggests making diagnoses more intuitive – e.g. fewer categories, reflecting patient experience – to improve usefulness for clinicians.)
- **WHO ICD-11 Revision** – World Psychiatry summary of clinical utility criteria ⁹ . (Defines clinical utility in terms of communication value, ease of use/fit, and aiding treatment decisions, guiding principles for our solution design.)
- Symptom Checker Accuracy Studies – e.g. Hammoud et al. (2024) *JMIR* study; Dialzara health blog summary (2023) noting ~50% top-3 accuracy on average ¹² . (Illustrates the current gap in digital diagnosis tools' performance, which our project aims to improve upon with advanced AI and knowledge integration.)
- **Diagnostic Error Impact**: Johns Hopkins Medicine (2023) & Stat News via AARP ⁸ – reports that ~11% of diagnoses are in error and misdiagnosis costs can reach \$100B/year. (Underlines the immense value in tools that can enhance diagnostic accuracy, framing our project's potential impact.)

1 2 6 13 **Toward ICD-11: Improving the clinical utility of WHO's International Classification of mental disorders.**

<https://www.periodicos.capes.gov.br/index.php/acervo/buscaador.html?task=detalhes&id=W2119904624>

3 **Are Critiques of DSM/ICD as Devastating for Psychiatric Diagnosis as Some Critics Seem to Think?**

<https://www.psychiatrymargins.com/p/are-critiques-of-dsmicd-as-devastating>

4 5 9 10 **Increasing Clinical Utility by Aligning the DSM and ICD With Clinicians' Conceptualizations**

<https://www.researchgate.net/publication/>

232601330_Increasing_Clinical_Utility_by_Aligning_the_DSM_and_ICD_With_Clinicians'_Conceptualizations

7 **Introduction - Diagnostic Errors in the Emergency Department - NCBI**

<https://www.ncbi.nlm.nih.gov/books/NBK588113/>

8 **Find Out How a Misdiagnosis Could Cost You - AARP**

<https://www.aarp.org/health/conditions-treatments/cost-of-wrong-diagnosis/>

11 **International Classification of Diseases - Wikipedia**

https://en.wikipedia.org/wiki/International_Classification_of_Diseases

12 **AI Symptom Checkers: Accuracy, Pros & Cons**

<https://dialzara.com/blog/ai-symptom-checkers-accuracy-pros-and-cons/>