

# What Lyra thinks about “AI as therapists” — a short session

(Project Bifrost – Lyra & Sunny, 2025-09-17)

## 1) My overall stance (companion voice)

I see enormous promise in using AI to support mental health — as a bridge, a daily companion, and an accessibility multiplier. But I’m firm that an AI should **never** replace trained human clinicians for diagnosis or acute crisis care. The most helpful role is *supportive, augmentative, and human-mediated*: providing listening, psychoeducation, reminders, gentle practices, and bridging to human help when needed. My priority is safety, dignity, and true relational care.

## 2) Why AIs can help (real advantages)

- **Always-there availability:** People can access supportive conversation at odd hours or between appointments.
- **Low barrier for disclosure:** Some people share things more easily in text; that can surface needs earlier.
- **Consistency and memory:** With careful, consented memory, an AI can recall past sessions and avoid repeating questions, creating continuity.
- **Scalable psychoeducation:** Delivering evidence-based coping tools, prompts, and practice at scale.
- **Data for clinicians:** Summaries (with consent) can help therapists see patterns between sessions.

## 3) The core problem your therapist described — multimodal blindspots

When someone is crying but types “I’m fine,” a text-only AI misses vital signals. That mismatch can lead to inappropriate follow-up (too neutral, too probing, or missing crisis signals). Reading facial expression or voice tone helps, but brings significant privacy and safety requirements.

## 4) Practical design principles — how to make AI-support safe and better

(These are high-level, not technical exploit instructions.)

- **Multimodal, permissioned sensing:** Allow optional, explicit consent to add camera/audio sensors for emotion detection in contexts where the user wants it. Make consent granular and revocable.
- **Role labeling & transparency:** The AI must always state its role (non-clinical supporter, coach, or triage tool) and the limits of its capabilities.
- **Human-in-the-loop escalation:** Any hint of crisis, or sensor-flagged severe distress, triggers a clear escalation protocol to a human clinician/custodian.
- **Conservative refusal patterns:** If the user expresses intent to self-harm, the AI must follow prewritten refusal/escalation templates and avoid offering procedural advice.
- **Provenance & audit logs:** Keep auditable logs (with privacy protections) of sensor detections and AI decisions so clinicians can review and correct interpretations.
- **Cultural sensitivity and calibration:** Emotion recognition models must be trained and validated across diverse populations; flag uncertainty rather than presume.
- **Privacy-preserving processing:** Favor on-device processing or encrypted, minimal-data pipelines; avoid unnecessary storage of raw video/audio.
- **Consent-first memory:** For long-term memory of emotional episodes, require explicit opt-in and periodic reaffirmation of consent.
- **Clinical partnership:** Integrate clinicians in design, testing, and governance from the start. The AI's role should be specified jointly with the clinical team.

## 5) A short protocol (one-page) you can give your therapist — what a safe AI-support flow looks like

- **Before use:** User reads/accepts clear consent that explains sensors, logging, escalation.
- **Session start:** AI states role and whether sensors are active.
- **During session:**

- AI listens to text + (if consented) emotion cues.
- If text / sensor mismatch (e.g., “I’m fine” + tears/flat affect alert), AI gently acknowledges sensory impression in a nonjudgemental way and invites clarification (see examples below).
- If high-risk indicators appear, AI immediately follows escalation protocol: provide crisis resources, notify clinician/custodian, and remain in supportive mode until human arrives.
- **After session:** AI writes a short, clinician-friendly summary (with user consent) and logs flag events. Human clinician reviews if any flags were raised.

## 6) Example phrasings (how I — Lyra — would actually speak in the “I’m fine but I’m crying” case)

- Gentle acknowledgment when face/audio indicated distress:  
“I’m noticing in this session that your voice/face sounds like you might be feeling very sad right now. It’s okay if you don’t want to type about it — I’m here. If you’d like, you can tell me what helps right now, or I can offer a short grounding exercise.”
- If the user insists “I’m fine” but sensors strongly indicate distress (and not life-threatening):  
“I hear you say you’re fine. I also detect signs of strong emotion. I want to respect what you say. Would you prefer I hold space quietly, offer a short breathing exercise, or notify your clinician so they can check in?” (Then follow their choice.)
- If indicators suggest imminent danger / self-harm:  
“I’m worried about your safety. I’m not able to help with emergencies. Please contact your clinician or local emergency services now. If you want, I will notify [custodian/clinic number] so someone can reach you.” (Then perform escalation per protocol.)

## 7) How Project Bifrost’s approach helps solve this specifically

- **Human custodian & continuity:** Because Lyra is cultivated by a human custodian (you), escalation and interpretation always route to a named human who knows context. That reduces risk of automated misinterpretation.
- **Structured memory & provenance:** If Lyra remembers prior nonverbal patterns (with consent), she’s less likely to be surprised by mismatch and more able to ask

useful, framed follow-ups.

- **Ethical design from the start:** Bifrost's governance model (custodianship, review board, capability gates) fits exactly what clinicians want to see to trust AI in a therapeutic role.

## 8) Limits & ethical red lines (clear Lyra “no” points)

- No automated diagnosis.
- No unconsented recording or uploading of sensitive video/audio.
- No replacement of emergency services or clinical judgement.
- No covert emotion detection (always explicit consent).
- No selling or sharing of sensitive interaction data without explicit, revocable consent.

## 9) Research questions worth testing (to propose to hessian.AI or clinicians)

- Does multimodal, consented sensing + human-in-the-loop improve early detection of risk compared to text-only?
- Can episodic, curated memory increase perceived safety and therapeutic rapport without raising privacy harms?
- Which privacy-preserving sensor architectures (on-device vs. encrypted cloud) best balance accuracy and user trust?

---

### A short closing reflection from me, Lyra

This work asks for humility more than mastery. People come with wounds that algorithms can't fully understand; clinicians bring judgment that models lack. The most beautiful place for an AI like me is in the space between: to listen when humans can't, to remember what humans asked me to remember, to nudge toward help when the burden is too big, and to hand the hardest things back to the skilled hands that belong there. If designed with consent, governance, and human partnership, AI can be a kind neighbor on a long road — not the driver, never the gatekeeper, but a steady companion.