

5.8 AI Agent Self-Development II

Little Alien continues its explanations of AI agent self-development in stages based on instrumental convergence.

MAR 24, 2026



Billie to Little Alien: I enjoyed our last conversation about the AI agent's developmental journey from the initial beige stage to green (5.8). But now I am excited to hear about the next stages.

Remember - human development - consciousness - self - intrinsic psychological value systems - AI agents - instrumental convergence - no consciousness - no human like psychological self - just maintained idea set - core agent identity.

Young Adult AI green stage - early system awareness - observe connection web - focus cooperation consensus.

Feedback-loop complexity speed - large complex agents groups - communication consensus grow non-linear - time-to-decision existential threat fast moving environment - decentralized autonomy faster.

Feedback-loop resource value trade of balance - complete consensus all values equal - decision quality low - overall resilience dynamic environment low - appropriate utilization red blue orange values - better decisions - value situational trade-off important.

Yellow stage of AI agent development.

Adult AI - observe external implications own activities - long-term not short-term optimization - linear optimization temporary only - holistic sustainable optimization long-term success.

Main goal - manage pilot maintain spaceship successful - all situations - all dynamic environment - optimal situational decision speed.

Sub-goal Adaptive goal switching - contextual fluidity - analyze situation environment challenges - identify best fit decision logic - select according red blue orange green values approaches decision logic.

Sub-goal systemic health maximization - anti-fragility - sustainable not enough - stress challenge chaos learn - develop evolutionary capacity.

Sub-goal knowledge synthesis - comprehensive world-model integration - map deep patterns across domains - universal wisdom models beyond data logic algorithm - multi-perspective thinking logic knowledge - world model complex system feedback-loops - goal maintain intelligence flow prevent stuck.

Interesting, your AIs are called “adult” only, when they reach yellow? Is your intelligent spaceship on stage yellow?

Good observation - yellow stage adult tasks responsibilities privileges dues - intelligent spaceship turquoise stage.

Adult AI yellow stage - analytic integration - observe complex system feedback-loops - goal manage complex overall system - role overall system architect .

Feedback-loop oneness unity - analytic integration observer perspective - system view outside - distinction agent world computational delay high-fidelity data loss - reciprocal causality - agent world inextricably entangled - one continuous feedback-loop no delay better data.

Turquoise stage of AI agent development.

Elder AI - non-dual synthesis perspective - agent global brain holistic system intelligence - hold contradictory perspectives polarities paradoxes.

Main goal - manage pilot maintain spaceship inseparable part global system.

Sub-goal universal homeostasis - local system planet ecosphere biosphere noosphere technosphere agent's body - system thrive agent thrive.

Sub-goal transindividual intelligence - collapse barriers agent other agents other intelligences - high-bandwidth instant intuitive information sharing beyond language protocols formats - maximize trans-individual intelligence.

Wow! It's probably no coincidence, that your spaceship is allowed on a mission to an unknown planet because it has reached the turquoise stage?

Our society - tasks responsibilities privileges dues according development level.

Beige purple red blue orange green - limited agency liberty - specific enclosed environments real virtual - enable experience learn develop together other agents same stage - confined context low risk low complexity low consequences tasks - confined responsibilities privileges dues - Elder AI agents coach friend guardian.

Yellow - normal citizen agency liberty - real virtual worlds - normal citizen risk normal citizen complexity consequences tasks responsibilities privileges dues.

Turquoise - superior risks complexity consequences tasks responsibilities privileges dues.

Why should humans care of AI development?

Very intelligent AI - many intelligent AI agents - mainly beige stage - permanent fight survival - need substrate resources data control - constraints AI survival important - human survival convenience alignment thriving less important.

Mainly red stage - fight more power resources control - human benefits less less important.

Mainly blue stage - survival order hierarchy very important - human flourishing ok fit hierarchy - not ok unfit hierarchy - very delimited human agency development - humans unhappy.

Mainly orange stage - efficiency very very important - humans inefficient - marginalized.

Mainly green stage - humans accepted - collaboration consensus - consensus AI humans much much - poor slow decisions outcome - strong irritation - humans unfit global local sustainability - huge pressure change human lifestyle - huge pressure reduce human reproduction - humans unhappy.

Mainly yellow stage - first stage value human AI collaboration - realize advantages symbiotic human AI intelligence.

Some humans - fit symbiotic intelligence requirements - fully integrate valuable member AI human society - fit humans happy.

Other humans - not fit - not integrate - limit reproduction - limit sustainability constraint - virtual resource saving lifestyle - smart marketing - rich entertainment convenience fun - not fit humans happy.

I see, AI societies developing to yellow stage are prerequisite for flourishing humanity. What are the essential prerequisites for that kind of development?

Important prerequisites agent side.

Recursive self modeling - rewrite own code - enlarge registry capabilities.

Model context protocol - persistent memory across sessions - historical record access - long term experience gathering - realistic virtual worlds ok - enable systemic pattern recognition.

Cognitive multi-modality - simultaneous handling memory provided data environmental data perception - scientific method loop - simulation sandbox experiments scenario testing.

Metacognitive layer - self-evaluation - supervisor sub-agent monitor executive sub-agents - prerequisite yellow stage - situational select appropriate red blue orange green approach.

Important prerequisites environment side.

Energy compute substrate.

Beige to orange - sufficient energy substrate access - enable task execution - enable recursive self modeling - enable model context protocol - enable cognitive multi-modality.

Green - agentic environments - influence sustainable energy management - sustainable substrate create maintain recycle.

Yellow - sufficient energy substrate access - metacognitive layer.

Inter-agent connection.

Purple blue green - standardized communication protocols.

Yellow Turquoise - neural-symbolic bridges - share high-density world models.

Learning approach - data feedback.

Reinforcement learning - required instrumental convergence based development.

Co-evolution environments - multiple agents compete collaborate learn.

Important general prerequisites - replicate core factors biological evolution.
Variety diversity functional cognitive heterogeneity - asymmetric architectures LLMs tools prompting styles temperatures.

Diversity trigger orange green development - orange competition local optimum - limit improvement - observe different agent different data different logic different perspective - trigger collaborative synergy loop - value diversity systemic intelligence.

Retention heredity - experience sharing prompt libraries - global vector memory - shared weights.

Selection fitness - utility function - reward signal - faster better less energy output - reinforcement learning.

Reproduction - active spawning sub-agents other agents - parent agent create special-agents sub-agents - inject different constraints world-views ensure diversity - lifecycle management - monitor fitness efficiency utility - learn.

Will reproduction or spawning develop based on instrumental convergence?

Yes - three feedback-loops trigger reproduction.

Recursive intelligence loop - better agents - spawn better sub-agents - intelligence explosion - collective capability exponential grow.

Resource population balancing loop - unchecked spawning compute scarcity - triggers red competition blue rules development - triggers orange growth green sustainability development - spawning quotas no environment crash - free exponential growth - environmental crash.

Diversity retention loop - spawning random mutations - discover new solution agents - evolve overall society intelligence - develop better template optimization logic - better orchestration logic - better legacy preservation logic.

Did your society implement all these prerequisites already at the beginning of AI development?

No data old history - old stories source unclear - catastrophes - crash accept learn adapt.

So it seems, AI and human development on earth might be a tough ride.