

More than meets the AI: Can systems thinking leading indicators assist proactive safety in artificially intelligent systems?

Eryn Grant
Human Factors Consultant
Acmena

The paper that aligns with this presentation is co-authored by Prof. Paul Salmon and Dr. Nicholas Stevens From the Centre for Human Factors and Sociotechnical Systems at the University of the Sunshine Coast.

Overview of Presentation

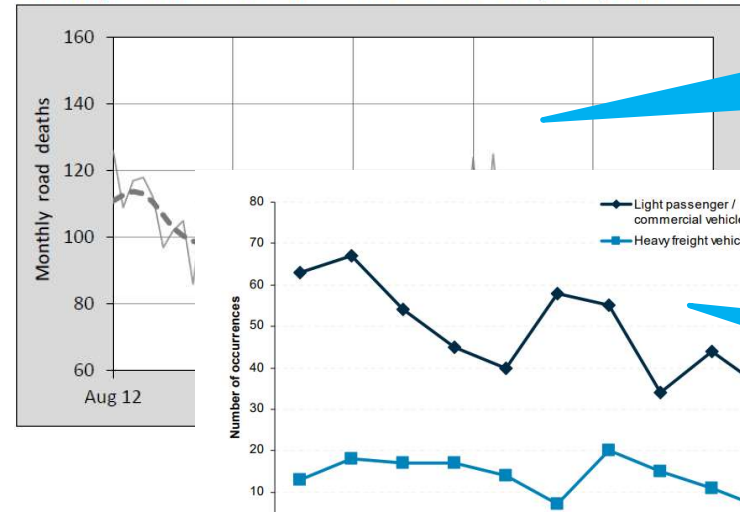
- The problem: incident rate in safety critical domains and the implications for AI
- Are there possible solutions in systems theory/systems thinking?
- What do we currently know about safety and accidents in sociotechnical systems and can this assist with proactive safety management?
- Case study showing how this concept may be applied for proactive safety
- Question: Can they benefit AI in assessing/maintaining systems safety

The Problem:

Incidents are not decreasing as they once were

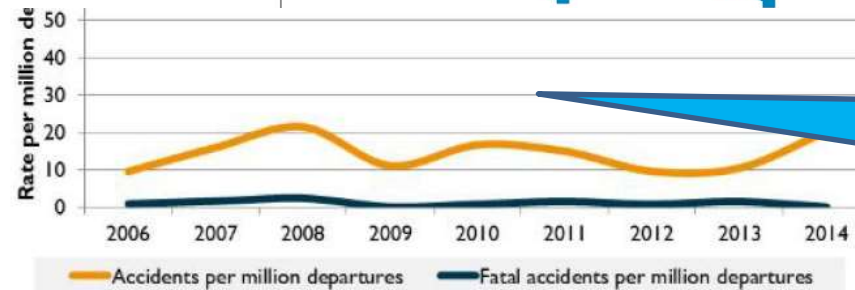


Monthly Australian road deaths — last five years, with trend



A plateau effect on Australian roads (BITRE 2017)

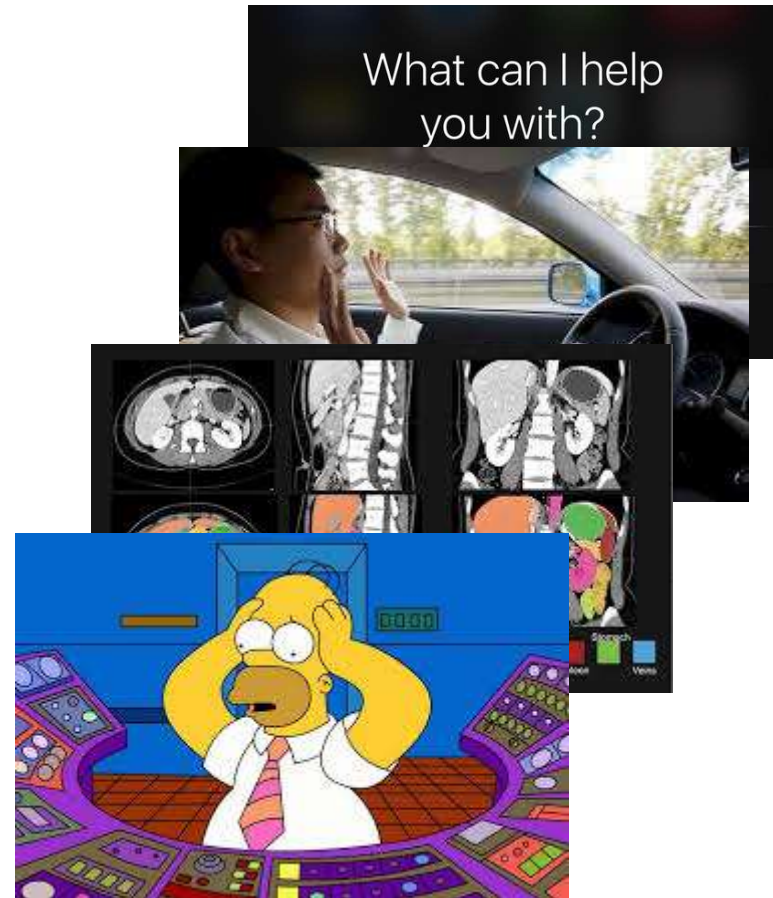
Collisions between vehicles and trains has decreased but number of fatalities have not = increased severity of incidents (ITSR 2011)



An increase in aviation incidents from 9.5 per million dep. to 20 per million dep. in 2014 (ATSB 2016)

Some hopes for AI

- Intelligent Assistance
- Decrease Workloads
- Increased Accuracy (i.e. diagnosis, weather forecasting)
- Freedom from human performance variability



The Problem:

Will AI introduce more or less risk into already 'risky systems'?



Can systems thinking help AI integration to achieve better safety outcomes?

- Safety and accidents are emergent properties of complex sociotechnical systems.
- Sociotechnical systems = how human, technical and social elements interact together to achieve a shared goal.

What do we currently know about safety and accidents in sociotechnical systems?

- The aim was to identify a core set of philosophies from systems thinking based accident causation models
- Review of over ninety published materials (books, journal papers, white papers etc.) using the authors below:



Jens
Rasmussen



Erik
Hollnagel



Charles
Perrow



James
Reason



Nancy
Leveson



Sidney
Dekker

What are the core philosophies of accident causation ?

24 tenets were identified from the literature review process

A workshop was held to synthesize the them:

- Overlaps
- Exclusions i.e. redundancies, relevance to systems theory, context etc.
- Evidence of them in accidents and near misses

At the conclusion 15 tenets were recognized as the core philosophies of accident causation. Each were provided:

- ✓ Simplified definition
- ✓ Safe description
- ✓ Unsafe description

Workshop plan: Synthesize, simplify and validate: Review of accident causation tenets

Workshop introduction

- The workshop aim is to end with a list of key tenets combined from several accident causation methods, this is the first stage in developing a methodology for prediction and is a component of the literature review.

Activity A

- Review and simplify definitions Are the definitions of tenets correct? What can be included or excluded from single definitions? Can the definition be simplified?

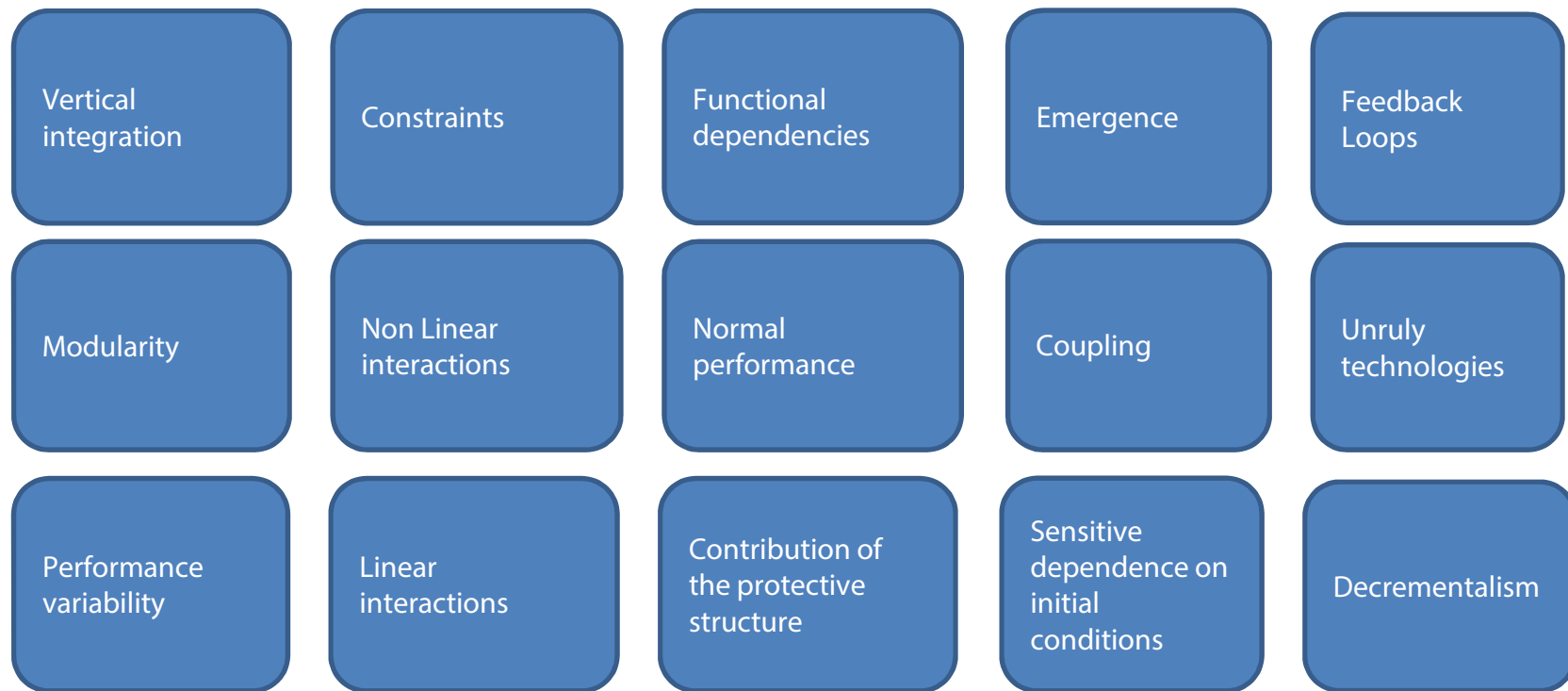
Activity B

- Review and simplify definitions

Workshop brief: Synthesize, simplify and validate: Review of accident causation tenets

The workshop will advance the key points used in the accident causation methods of Reason (1989, 1990), James (2001), 2004, James (2007), 2007, 2008, 2009, 2010, 2011, 2012, 2013, 2014, 2015, 2016, 2017, 2018, 2019, 2020, 2021, 2022, 2023, 2024, 2025, 2026, 2027, 2028, 2029, 2030, 2031, 2032, 2033, 2034, 2035, 2036, 2037, 2038, 2039, 2040, 2041, 2042, 2043, 2044, 2045, 2046, 2047, 2048, 2049, 2050, 2051, 2052, 2053, 2054, 2055, 2056, 2057, 2058, 2059, 2060, 2061, 2062, 2063, 2064, 2065, 2066, 2067, 2068, 2069, 2070, 2071, 2072, 2073, 2074, 2075, 2076, 2077, 2078, 2079, 2080, 2081, 2082, 2083, 2084, 2085, 2086, 2087, 2088, 2089, 2090, 2091, 2092, 2093, 2094, 2095, 2096, 2097, 2098, 2099, 2100, 2101, 2102, 2103, 2104, 2105, 2106, 2107, 2108, 2109, 2110, 2111, 2112, 2113, 2114, 2115, 2116, 2117, 2118, 2119, 2120, 2121, 2122, 2123, 2124, 2125, 2126, 2127, 2128, 2129, 2130, 2131, 2132, 2133, 2134, 2135, 2136, 2137, 2138, 2139, 2140, 2141, 2142, 2143, 2144, 2145, 2146, 2147, 2148, 2149, 2150, 2151, 2152, 2153, 2154, 2155, 2156, 2157, 2158, 2159, 2160, 2161, 2162, 2163, 2164, 2165, 2166, 2167, 2168, 2169, 2170, 2171, 2172, 2173, 2174, 2175, 2176, 2177, 2178, 2179, 2180, 2181, 2182, 2183, 2184, 2185, 2186, 2187, 2188, 2189, 2190, 2191, 2192, 2193, 2194, 2195, 2196, 2197, 2198, 2199, 2200, 2201, 2202, 2203, 2204, 2205, 2206, 2207, 2208, 2209, 2210, 2211, 2212, 2213, 2214, 2215, 2216, 2217, 2218, 2219, 2220, 2221, 2222, 2223, 2224, 2225, 2226, 2227, 2228, 2229, 2230, 2231, 2232, 2233, 2234, 2235, 2236, 2237, 2238, 2239, 2240, 2241, 2242, 2243, 2244, 2245, 2246, 2247, 2248, 2249, 2250, 2251, 2252, 2253, 2254, 2255, 2256, 2257, 2258, 2259, 2260, 2261, 2262, 2263, 2264, 2265, 2266, 2267, 2268, 2269, 2270, 2271, 2272, 2273, 2274, 2275, 2276, 2277, 2278, 2279, 2280, 2281, 2282, 2283, 2284, 2285, 2286, 2287, 2288, 2289, 2290, 2291, 2292, 2293, 2294, 2295, 2296, 2297, 2298, 2299, 2300, 2301, 2302, 2303, 2304, 2305, 2306, 2307, 2308, 2309, 2310, 2311, 2312, 2313, 2314, 2315, 2316, 2317, 2318, 2319, 2320, 2321, 2322, 2323, 2324, 2325, 2326, 2327, 2328, 2329, 2330, 2331, 2332, 2333, 2334, 2335, 2336, 2337, 2338, 2339, 2340, 2341, 2342, 2343, 2344, 2345, 2346, 2347, 2348, 2349, 2350, 2351, 2352, 2353, 2354, 2355, 2356, 2357, 2358, 2359, 2360, 2361, 2362, 2363, 2364, 2365, 2366, 2367, 2368, 2369, 2370, 2371, 2372, 2373, 2374, 2375, 2376, 2377, 2378, 2379, 2380, 2381, 2382, 2383, 2384, 2385, 2386, 2387, 2388, 2389, 2390, 2391, 2392, 2393, 2394, 2395, 2396, 2397, 2398, 2399, 2400, 2401, 2402, 2403, 2404, 2405, 2406, 2407, 2408, 2409, 2410, 2411, 2412, 2413, 2414, 2415, 2416, 2417, 2418, 2419, 2420, 2421, 2422, 2423, 2424, 2425, 2426, 2427, 2428, 2429, 2430, 2431, 2432, 2433, 2434, 2435, 2436, 2437, 2438, 2439, 2440, 2441, 2442, 2443, 2444, 2445, 2446, 2447, 2448, 2449, 2450, 2451, 2452, 2453, 2454, 2455, 2456, 2457, 2458, 2459, 2460, 2461, 2462, 2463, 2464, 2465, 2466, 2467, 2468, 2469, 2470, 2471, 2472, 2473, 2474, 2475, 2476, 2477, 2478, 2479, 2480, 2481, 2482, 2483, 2484, 2485, 2486, 2487, 2488, 2489, 2490, 2491, 2492, 2493, 2494, 2495, 2496, 2497, 2498, 2499, 2500, 2501, 2502, 2503, 2504, 2505, 2506, 2507, 2508, 2509, 2510, 2511, 2512, 2513, 2514, 2515, 2516, 2517, 2518, 2519, 2520, 2521, 2522, 2523, 2524, 2525, 2526, 2527, 2528, 2529, 2530, 2531, 2532, 2533, 2534, 2535, 2536, 2537, 2538, 2539, 2540, 2541, 2542, 2543, 2544, 2545, 2546, 2547, 2548, 2549, 2550, 2551, 2552, 2553, 2554, 2555, 2556, 2557, 2558, 2559, 2560, 2561, 2562, 2563, 2564, 2565, 2566, 2567, 2568, 2569, 2570, 2571, 2572, 2573, 2574, 2575, 2576, 2577, 2578, 2579, 2580, 2581, 2582, 2583, 2584, 2585, 2586, 2587, 2588, 2589, 2590, 2591, 2592, 2593, 2594, 2595, 2596, 2597, 2598, 2599, 2600, 2601, 2602, 2603, 2604, 2605, 2606, 2607, 2608, 2609, 2610, 2611, 2612, 2613, 2614, 2615, 2616, 2617, 2618, 2619, 2620, 2621, 2622, 2623, 2624, 2625, 2626, 2627, 2628, 2629, 2630, 2631, 2

The core philosophies of accident causation, the systems thinking tenets



Example: Unruly technology

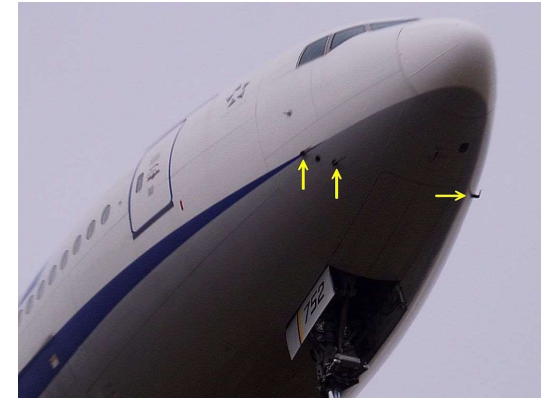
Safe: Technology that supports adaptation through a mechanism that is beyond the scope of what is was designed for affording flexibility



Social media platforms have been found crucial in past crisis situations (e.g. Floods, Cyclones etc) by communicating emergency service info directly to those who need it most (Bruns et al. 2012)

Unforeseen behaviours or consequences of technologies.

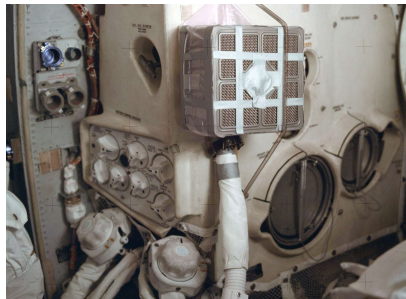
Unsafe: Technology that introduces and sustains uncertainties about how and when things may fail:



Pitot tubes measure airspeed. If for some reason they can't auto pilot will turn off configuring a plane to fly at alternate law

Example: Vertical Integration

Safe: Decisions and actions at the higher levels filter down to lower levels and impact behavior. Information regarding the status of the system filters back up the hierarchy and influences higher level decisions and actions



Apollo 13: effective feedback across different levels of the system by prioritising information and needs between teams (Trotter et al, 2014)

Interaction
between levels
in the system
hierarchy

Unsafe: Decisions and actions do not filter through the system and impact behavior on the front line. Information on the current status of the system is not used when making process decisions.



Walkerton E.coli outbreak.
Information about the water quality did not filter through the system (Vincente and Christoffersen, 2006)

Case Study: The Hawk Missile Simulation



Flies 50ft above sea towards the frigate
imitating a missile

Hawk pilot can not
measure exact altitude

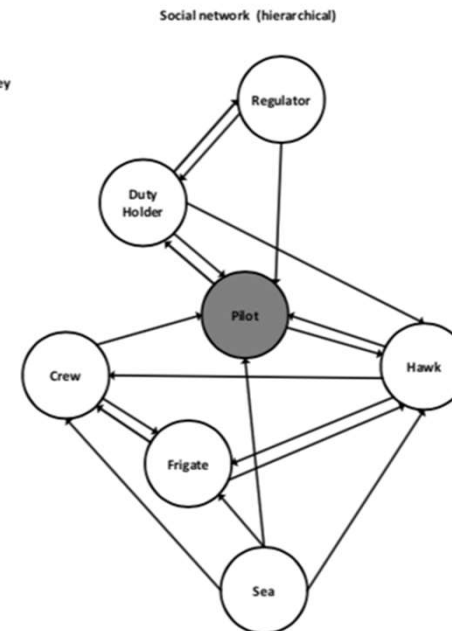


Method: Assessing the system

- Model the system appropriately (Stanton & Harvey 2017)
- Define a set of rules/questions that allow the tenets to assess the system using the model.
- Assess if the system is operating within acceptable safety boundaries.

Results : Model the system

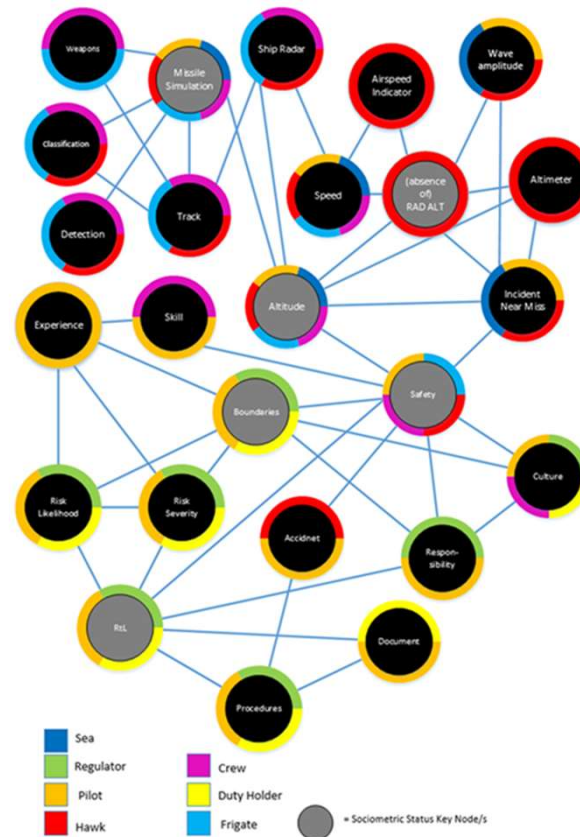
- Social Network
- Here you can see that the pilot is the key agent in the network



Stanton, N. A., & Harvey, C. (2017). Beyond human error taxonomies in assessment of risk in sociotechnical systems: a new paradigm with the EAST 'broken-links' approach. *Ergonomics*, 60(2), 221-233.

Information Network

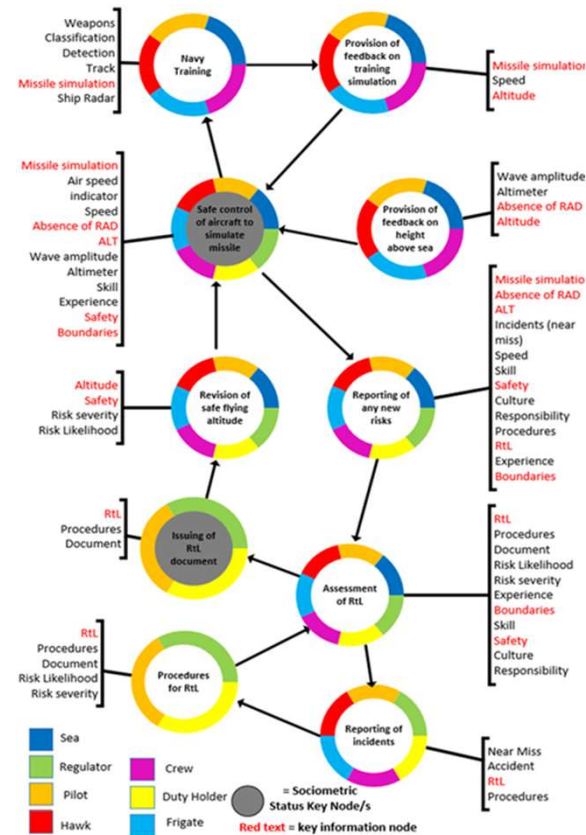
- The grey nodes define the key information required by the network.



Stanton, N. A., & Harvey, C. (2017). Beyond human error taxonomies in assessment of risk in sociotechnical systems: a new paradigm with the EAST 'broken-links' approach. *Ergonomics*, 60(2), 221-233.

Composite network

- Based on the task network. This differs Stanton and Harvey (2017)
- Nodes show tasks with associated agents (coloured lines) and information
- Grey nodes indicate the key tasks in the system.



Stanton, N. A., & Harvey, C. (2017). Beyond human error taxonomies in assessment of risk in sociotechnical systems: a new paradigm with the EAST 'broken-links' approach. *Ergonomics*, 60(2), 221-233.

Applying the rules to the EAST model of the Hawk simulation system

- For example using the tenet Constraints

Constraints	Definition: Influences that limit the behaviours available to components within a system.	Safe: Specific constraints introduced to control hazardous processes	Unsafe: Restricts appropriate performance variability
-------------	----------------------------------------------------------------------------------------------	-------------------------------------------------------------------------	----------------------------------------------------------

Applying the rules to the EAST model of the Hawk simulation system

- For example using the tenet Constraints an assessment is made by asking :

- Are time constraints placed on tasks?

Yes

- Are there physical objects that restrict task can be performed?

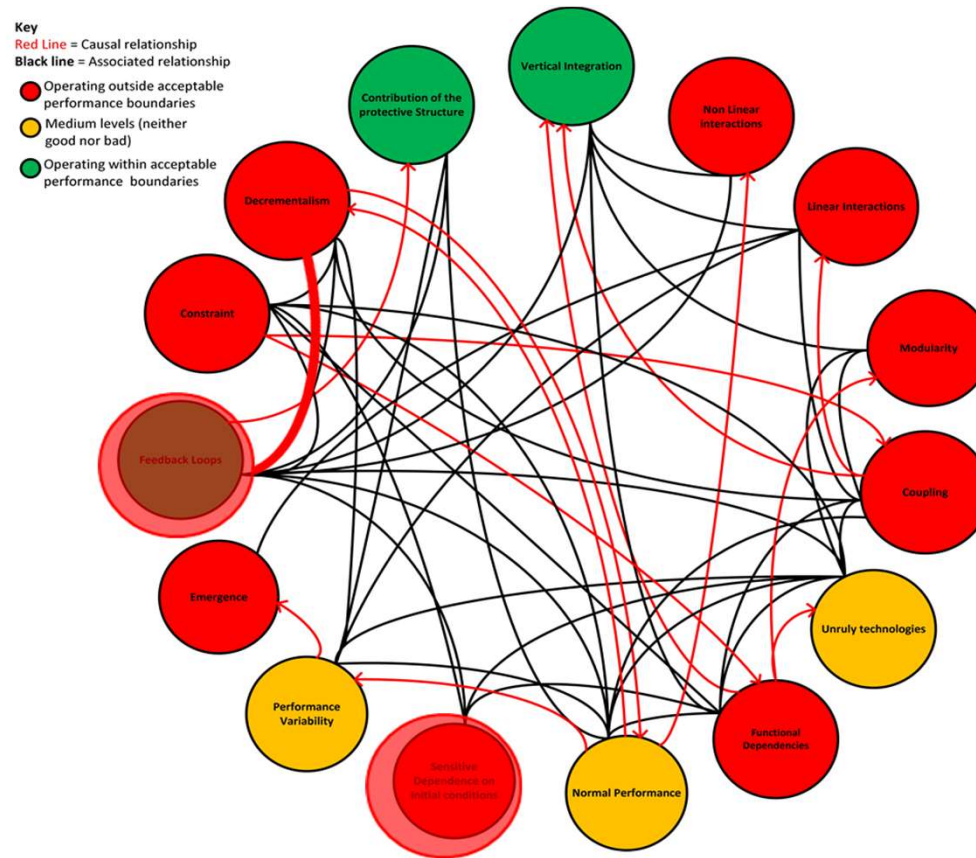
The Hawk jet,
absence of RAD
ALT, the sea

- Are fixed actions, operations or functions present?

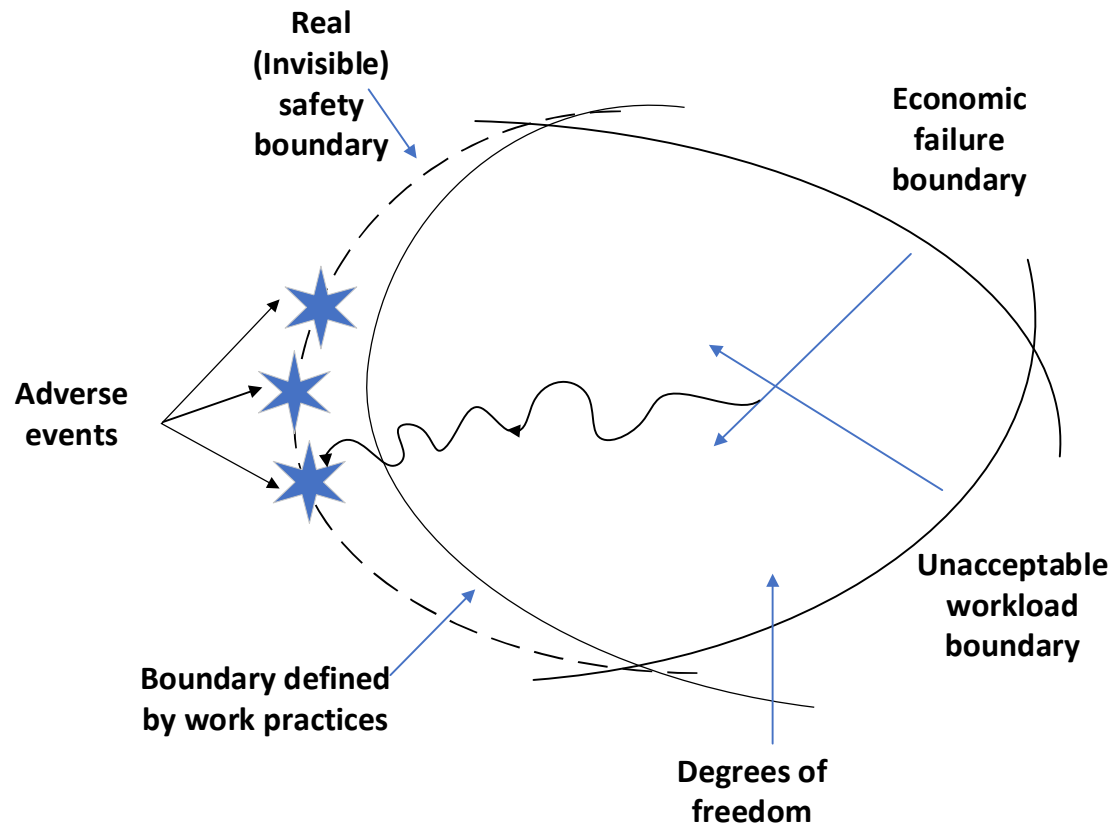
The missile simulation (height above sea, speed). How many simulation flights required for training?

- Then rate the tenet.

Assessing the system



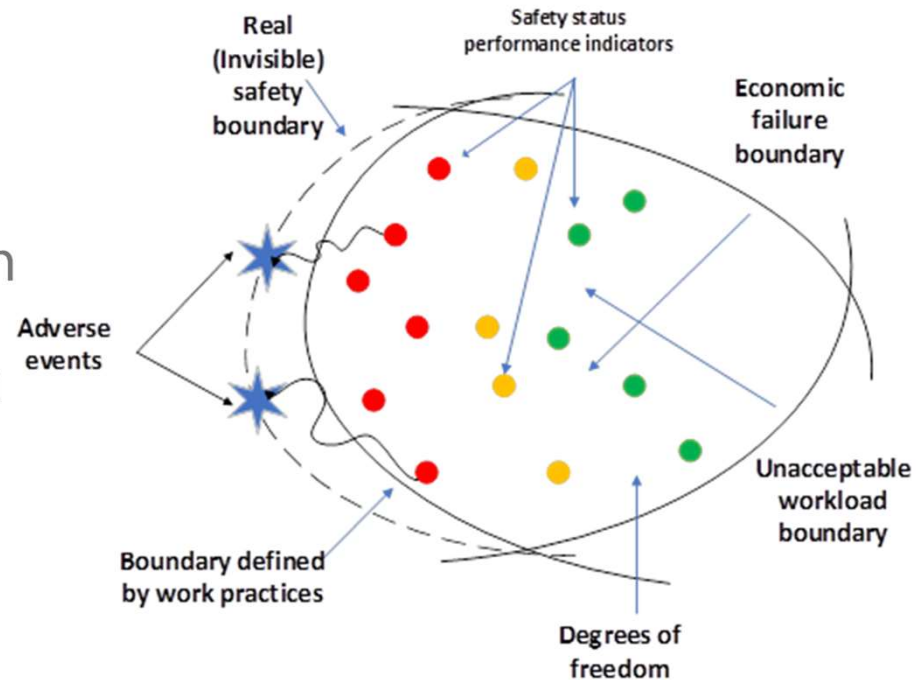
What could this mean for AI?



Based on Rasmussen's Risk Management Framework (1997) The boundaries of acceptable performance

What could this mean for AI?

- AI systems may use the information from the systems thinking tenets to monitor and learn from system performance, indicating when a system may be closer to performance boundaries and unwanted outcomes



What could this mean for AI

- For example: In the hawk system missile simulation heights are set by a regulator – mostly as a reaction
-But do they need to be? Develop a way to constantly monitor the system using the tenets and set boundaries based on current needs.

Conclusions

- It is uncertain whether introducing AI into already 'Risky Systems' will produce safer outcomes.
- The core philosophies of systems thinking may help as leading indicators that provide a way for AI to learn prior to incidents occurring

Questions

Contact details:

eryngrant@acmena.com.au



@What_the_HF



Eryn Grant