



MICKAITM

MICKAI EBOOK SERIES · No. 17

The Offline-First Doctrine.

Designing intelligence that keeps working when the link goes down.

AUTHOR

Micky Irons

Founder and named inventor, Mickai LTD.

19 June 2026 · v1 · mickai.co.uk

EBOOK · No. 17 IN A SERIES OF 34

Mickai LTD · Companies House 17166618 · press@mickai.co.uk · mickai.co.uk
UK IPO register, named inventor Mickle Wagstaff-Irons · Trade mark UK00004373277

TABLE OF CONTENTS

Contents

Foreword

A note from the author

Part One: The Dependency Trap

The Quiet Assumption Underneath Everything
Why Connectivity Became a Single Point of Failure
What We Lose When the Mind Lives Elsewhere

Part Two: The Mechanism of Graceful Degradation

Local-First Is an Architecture, Not a Fallback
Degrade Gracefully, Function by Function
Trust Without the Link: Sealing and Provenance

Part Three: Evidence, Economics, and the Hardware Question

The Economics of Owning Your Intelligence
Building Everything, Gated to the Hardware You Have
Provenance as Evidence You Can Defend

Part Four: Putting the Doctrine to Work

For the People Who Build Systems
For the People Who Buy and Deploy Them
For the People Who Write the Rules

Appendix

About the author

FOREWORD

A note from the author

I have spent the last few years building a system that assumes the network will fail. Not might fail. Will fail. That assumption sits at the centre of everything I design, and it puts me at odds with a large part of the industry. The prevailing model of intelligence today is a thin client talking to a distant data centre, a conversation that only works while the link holds. I think that is a fragile way to build anything that matters, and this short book is my attempt to explain why, in plain terms, with the engineering shown rather than asserted.

My name is Micky Irons. I am the founder of Mickai and the named inventor on the patent applications that protect it. Mickai is a Sovereign Intelligence Operating System, a full stack that runs fifty specialised models, twenty-five domain and twenty-five operational, on hardware the operator owns and controls. I did not arrive at offline-first as a marketing position. I arrived at it because I kept watching capable systems become useless the moment they were disconnected, and I decided that any intelligence worth trusting in a hard moment has to be able to think without asking permission from somewhere else.

This book is built as an argument rather than a catalogue. The first part sets out the problem with treating connectivity as a dependency. The second explains the mechanism, how a system can degrade gracefully instead of collapsing. The third examines the evidence and the economics, because a doctrine that costs more than it returns is just a preference. The fourth turns the argument into instructions you can act on, whether you build systems, buy them, or write the rules they have to follow.

A word on tone before we begin. I am not interested in frightening anyone, and I am careful not to overclaim. Where a capability is designed and filed but not yet running in production, I say so. Where a standard belongs to someone else, I credit them. What I want to leave you with is not alarm but a different default, a habit of asking one question of every intelligent system you depend on. What happens to this when the link goes down? If the answer is nothing good, you do not have a tool. You have a hostage.

Micky Irons

Founder and named inventor, Mickai LTD · 19 June 2026

PART ONE: THE DEPENDENCY TRAP

Most intelligence today stops thinking the moment its connection drops, and that is a design choice, not a law of nature.

The Quiet Assumption Underneath Everything

Walk through almost any modern system that claims to be intelligent and you will find the same hidden clause in its design. It assumes the link is always there. The model lives in a data centre, the device in your hand is a window, and the conversation between them is the product. While the connection holds, the experience feels effortless. The assistant answers, the camera recognises, the document is summarised. None of it advertises the dependency, because the dependency is invisible until it is broken.

I call this the quiet assumption because nobody writes it down. No product page says this feature requires an uninterrupted route to a server two thousand kilometres away. Yet that is exactly what most of them require. The intelligence is not in the room with you. It is rented, streamed, and revocable, and the terms of that arrangement are written by the network, not by you. When the network is healthy you never read the terms. When it is not, you discover them all at once, usually at the worst possible moment, because the moments that break a connection are the same moments that strain everything else.

There is a reason the industry settled here, and it is not stupidity. Centralising the model makes it cheaper to update, easier to meter, and simpler to control. A single large model in a single place can serve millions, and every improvement reaches everyone at once. Those are real advantages and I do not dismiss them. The trouble is that they are advantages for the operator of the data centre, and the costs land on the person at the edge, the one whose link is thin, intermittent, or watched.

Intelligence you cannot reach is not intelligence you have, it is intelligence someone is willing to lend you for as long as the conditions suit them.

Consider where this falls apart. A clinic in a region with two hours of mains power a day. A vessel a hundred miles past the last cell tower. A factory floor wrapped in steel that swallows signal. A field team somewhere the link is not merely weak but actively contested, throttled or cut by someone who would rather they could not think clearly. In each case the centralised model is not slow. It is absent. And absence is a far worse failure mode than slowness, because there is nothing to wait for and nothing to fall back on.

Once you notice the quiet assumption you start seeing it everywhere, and you start resenting it. The smart speaker that goes deaf in a blackout. The translation app that cannot translate on the long-haul flight where you most need it. The driver assistance that degrades the instant the maps stop loading in a tunnel. None of these were designed to fail in these moments. They were designed for the happy path, and the happy path was defined by people in well-connected offices who had never lost their link for long enough to care.

Why Connectivity Became a Single Point of Failure

Engineers are trained to hunt for single points of failure and remove them. We add redundant power supplies, mirror databases, and spread services across regions so that no one machine can take the whole system down. Then, having done all that careful work, we route every act of intelligence through one fragile thread, the connection between the user and the cloud, and we call it acceptable because it usually works. We hardened everything except the one wire the whole product hangs from.

The link is fragile in more ways than people credit. It fails physically, when a cable is cut or a tower loses power. It fails economically, when data is metered and the budget runs dry mid-task. It fails politically, when an authority decides certain traffic will not pass. And it fails silently, in the degraded middle state where packets crawl and time out, which is often worse than a clean outage because the system keeps trying and never quite succeeds. A doctrine that only survives the clean cases is no doctrine at all.

What makes connectivity such a dangerous single point of failure is that it concentrates so much behind it. When intelligence lives in the cloud, the link does not just carry a request. It carries the model, the reasoning, the memory, and the authority to act, all in one channel. Cut that channel and you have not lost a feature. You have lost the mind of the system, its recall, and its right to do anything useful, simultaneously. The blast radius of a dropped connection is the entire product.

The contested edge is not an edge case

We tend to file disconnection under exceptions, the rare event we will handle later. I think this is backwards. For a large share of the world, and for almost everyone in a genuine emergency, degraded connectivity is the normal condition, not the exception. The edge that is contested, congested, or simply remote is where the most consequential decisions get made, and it is exactly where the centralised model is weakest. Designing for the data centre and patching for the edge gets the priorities the wrong way round.

If your system only works where the link is strong, you have built a tool for the places that already had the fewest problems.

None of this means the cloud is bad. It means the cloud is a place, and no single place should hold your ability to think. The error is not in using remote resources. The error is in depending on them for the core act of intelligence, so that their absence is your incapacity. Once you frame it that way, the fix becomes obvious in principle even if it is hard in practice. Move the thinking to where the thinker is,

and let the network add reach on top of an ability you already hold.



The Mickai pantheon.

What We Lose When the Mind Lives Elsewhere

There is a cost to renting your intelligence that goes beyond the moments it disappears. When the model lives elsewhere, so does everything you feed it. Your questions, your documents, your patterns of work, your mistakes, all of it travels to a place you do not control and lingers there under terms you did not write. The dependency is not only operational. It is a steady, quiet export of the very context that makes the intelligence valuable, and once that context has left the room you cannot call it back.

This is the part operators feel last and regret most. A system that needs the network to think also tends to need the network to remember, and a memory you do not hold is a memory that can be read, changed, or withdrawn. For a hospital, a court, a defence team, or a business with something worth protecting, that is not a convenience problem. It is a sovereignty problem, and sovereignty lost by default is the hardest kind to win back, because by the time you notice it is gone the leverage has already moved to someone else.

There is also a subtler loss, the loss of determinism. When intelligence is served from a remote system that updates on someone else's schedule, the thing that answered you yesterday may not be the thing that answers you today. The model can shift beneath you without notice. For casual use that is tolerable. For a system that has to be auditable, where you may one day have to explain exactly why a decision was made, an intelligence that quietly rewrites itself is a liability dressed as a feature.

I built Mickai around the opposite stance. The fifty specialised brains, twenty-five domain and twenty-five operational, run on the operator's own hardware on a substrate we call Poseidon, and they are built to keep working offline. The context stays in the room. The memory is held, not borrowed. When the system answers, it answers from a known state you can pin down and account

for later. That is not nostalgia for local software. It is a deliberate choice about where power should sit, made by someone who would rather own a smaller mind in the room than rent a larger one he cannot reach.

The question is never only can it think, it is also who can switch it off, and from how far away.

So this first part ends with a simple reframing. The problem is not that cloud intelligence is weak. It is often extraordinarily strong. The problem is that strength delivered through a single fragile channel becomes weakness at exactly the moments strength is needed most. The rest of this book is about how to keep the strength and remove the fragility, starting with the mechanism that makes graceful degradation possible rather than merely aspirational.

PART TWO: THE MECHANISM OF GRACEFUL DEGRADATION

A system can be built so that losing the link costs you capacity, not capability, and this part shows how.

Local-First Is an Architecture, Not a Fallback

The common way to handle disconnection is to bolt on a fallback. The real system lives in the cloud, and when the cloud is unreachable a shrunken offline mode steps in, usually a cache of yesterday's answers and an apology. This is better than nothing, but it betrays the underlying belief that offline is the degraded case. I want to invert that belief. In an offline-first architecture, local is the default and the network is the enhancement, not the other way around.

The distinction is more than philosophical. If you treat local as the fallback, you build the local path last, test it least, and let it rot. If you treat local as the default, the on-device intelligence is the thing you build first and trust most, and the network becomes an optional accelerant that adds reach when it happens to be there. The same words, local and cloud, describe both designs. The order of priority is what separates a resilient system from a fragile one wearing a resilience badge.

Offline-first does not mean offline-only, it means the link is allowed to improve the system but never to be the reason it works at all.

Concretely, this means the model that matters runs where the user is. In Mickai, the fifty brains are served from the operator's own machines, with a kernel of operational brains, the Chronus Kernel, scheduling and routing the work locally. Inference happens on hardware in the room. The reasoning, the recall, and the authority to act all sit there. When connectivity is present it can do useful things, fetch an update, synchronise a record, reach a colleague. When connectivity is absent, none of the core functions stop, because none of them were ever leaning on the link to begin with.

Building this way is harder, and I will not pretend otherwise. You have to fit capable models onto real hardware, manage memory and heat, and accept that the device in the room is not an infinite data centre. We do this by fine-tuning and specialising open foundations, models such as Llama 3.2 and Qwen 2.5, into focused brains that do their narrow job extremely well rather than one enormous model that does everything passably. Specialisation is what makes local-first practical, because a brain that only has to be excellent at one domain can be small enough to run where it is needed, and several small excellent brains beat one large mediocre one at the edge.

I want to be precise about the present state. We are actively training our own models now, building a sealed corpus and specialising those open foundations, and the funded roadmap scales that work toward fully native weights over time. That is the direction of travel, and I describe it as a direction rather than a finished destination on purpose. The architecture, though, does not wait for that journey to complete. Local-first is true today, with the models we run today, on the hardware operators own today.



The Mickai pantheon.

Degrade Gracefully, Function by Function

Graceful degradation is the heart of the doctrine, and it is worth defining carefully because the phrase gets used loosely. A system degrades gracefully when the loss of a resource removes a proportionate amount of capacity and nothing more. Lose the network and you should lose the things that genuinely need the network, while everything that does not need it carries on. The failure should be local to the function that depended on the missing resource, never global to the whole system.

Contrast that with what most systems actually do, which is collapse. One missing dependency cascades into total failure because everything was wired through the same thread. Graceful degradation is the engineering discipline of cutting those wires deliberately, so that each capability fails alone. It is unglamorous work. It means asking, for every feature, what does this truly require, and refusing to let it quietly acquire a dependency on the network that it does not actually need.

Tiers of capability

The practical pattern is to think in tiers. There is a core tier that must work with nothing but local hardware and power, the functions a user would be in real trouble without. Above that sits a tier that works better with a local network or peer connections but does not need the wider internet, so a clinic or a ship can still share between its own machines. Above that sits a tier that genuinely benefits from

full connectivity, the things that legitimately reach out to the world. Designing in tiers forces you to be honest about which features belong where, and most belong far lower than their builders assumed.

A well-built system does not ask whether it is online, it asks what it can do right now and does all of it.

There is a further discipline that matters here, which is honest signalling. When a capability is unavailable because the hardware or the link cannot support it, the system should say so plainly and tell the operator what would change that, rather than failing in a way that looks like a fault. We build Mickai so that an over-specified feature on a modest machine reports clearly that it needs more, with guidance, instead of pretending or silently breaking. Honesty about limits is part of resilience, because a user who knows the boundary can plan around it instead of being ambushed by it.

Done well, graceful degradation changes how disconnection feels. Instead of a wall, you meet a gentle slope. The system gets quieter, not dumber. The reach narrows, but the mind stays. A user who loses the link notices that some outward-facing things have paused, while the intelligence in front of them keeps reasoning, keeps remembering, and keeps acting on what it already has. That experience, capacity reduced but capability intact, is the entire point of the doctrine.

Trust Without the Link: Sealing and Provenance

A hard question sits underneath offline-first. If a system can act entirely on its own, disconnected and unsupervised, how do you trust what it did? When the cloud is removed, you also remove the central log that many systems lean on to prove what happened. An intelligence that thinks alone in the dark needs a way to make its actions accountable that does not itself depend on being online. Without that, sovereignty becomes a synonym for opacity, and that is not a trade I am willing to make.

Our answer is to seal every consequential action at the moment it happens, locally, into an Open Audit Record. The record is signed using a post-quantum digital signature under FIPS 204, the ML-DSA-65 scheme. I want to be clear that this is a NIST standard. We did not invent the cryptography. We adopted a published, vetted standard and built our auditing around it, which is exactly what a serious system should do rather than rolling its own scheme. The signing happens on the device, so the proof is created whether or not there is a network present.

An action you cannot prove is an action you cannot defend, so the proof has to be made at the edge, in the moment, with no link required.

Sealing locally solves the immediate problem of accountability, but it leaves a second question. How do you anchor those local records to something the wider world can check, when the wider world reconnects? For that we anchor provenance to Pantheon, a sovereign Layer 1 anchored to Bitcoin, structured as a base chain with fifteen application chains and a fixed-supply PAN token. The point of anchoring is not speculation. It is to give a local, offline-made record a path to a durable, independently

verifiable home once a link exists, without ever having needed that link to create the record in the first place.

This is the synthesis that makes offline-first trustworthy rather than merely independent. The system can think and act without the network, and yet every consequential thing it does is sealed with a standard post-quantum signature the instant it occurs, ready to be anchored when connectivity returns. Disconnection no longer means darkness. It means a stack of sealed records, accumulating quietly, each one provable on its own terms. The audit trail outlives the outage, and the order of events is fixed at the time they happened rather than reconstructed afterward.

There is a body of designed-and-filed capability that extends this further, custody mechanisms such as a dead-man's switch, key rotation, trustee succession, and post-quantum protection for the keys themselves. I describe these as part of the architecture because they are designed and filed, and I am careful not to present any pending feature as if it were already running in production. The principle behind them is consistent with everything in this part. Trust, like intelligence, should not depend on a permanent connection to anyone.



The Mickai pantheon.

PART THREE: EVIDENCE, ECONOMICS, AND THE HARDWARE QUESTION

Resilience has to pay for itself, and this part weighs the real costs against the real returns.

The Economics of Owning Your Intelligence

A doctrine that only makes engineering sense is a hobby. To be worth adopting, offline-first has to make economic sense as well, and the honest starting point is that it costs more up front. Owning hardware capable of local inference is a real expense. Specialising and serving your own models is a real effort. Anyone who tells you sovereignty is free is selling you something. So the question is not whether it costs, but whether what it buys is worth the price over the life of the system.

The rented model has a seductive cost curve at the start. Little capital outlay, pay as you go, scale on demand. What that curve hides is that you are renting in perpetuity, and the meter never stops. Every query, every token, every stored record carries a recurring charge, and the price is set by someone whose interests are not aligned with yours. Over years, a system used heavily can pay for its own hardware several times over and still owe rent next month. The cheap option is cheap the way a payday loan is cheap.

Renting intelligence feels affordable right up to the point where you realise you will be paying for it forever and own nothing at the end.

Owning your intelligence inverts that curve. The capital cost is real and lands early, but after it lands the marginal cost of thinking approaches the cost of electricity. There is no per-query meter, no surprise price change, no usage tier to be bumped into, and no vendor able to depreciate the model you depend on. For an operator whose work involves heavy, sustained, sensitive use, the crossover comes sooner than they expect, and after the crossover every additional unit of work is almost free. The economics reward exactly the users who need resilience most, the heavy and the serious.

There is also a cost that never appears on the rented invoice, the cost of the outage itself. When a centralised system goes dark, the loss is not the unpaid query. It is the decision not made, the patient not assessed, the operation paused, the opportunity gone. That cost is invisible in normal accounting because it shows up as absence, and absence does not generate a line item. Offline-first is, in part, a way of pricing in the outages that everyone else pretends will not happen, by making sure the work simply continues while the link is down.

Building Everything, Gated to the Hardware You Have

The obvious objection to local-first is hardware. The cloud offers, in principle, unlimited capacity, and the machine in the room does not. How can a system that runs on owned hardware ever keep up with one that draws on a data centre the size of a town? It is a fair objection, and the answer is not to pretend the constraint away. The answer is to design for the full range of hardware and let each capability scale to what is actually present.

Our principle is to build every feature, even ones that exceed what a given machine can run today, and to gate them to the hardware in front of them. On a modest workstation, the demanding features report that they need more capacity and explain what an upgrade would unlock, rather than being silently absent. The same software, moved onto more capable hardware, simply does more. The capability is always present in the design. What varies is how much of it the current silicon can serve, and the system is honest about that boundary at all times.

A lineup, not a single box

This is why we think in terms of a hardware lineup rather than a single device. At the modest end, a single workstation with a 24GB professional GPU runs a meaningful set of the brains locally and serves a real day's work offline. At the flagship end sits a far larger server class machine, the Prometheus Edge Server, with multiple Blackwell-class accelerators and multi-terabyte memory, the kind of system that can hold many of the fifty brains resident at once and run the heaviest functions without reaching for the cloud at all. The doctrine does not demand the flagship. It demands that the architecture scale smoothly between the two, so an operator buys the resilience their work justifies.

The right amount of hardware is the amount that lets you keep working through the outage you are actually likely to face.

Specialisation is what makes this scaling realistic. A single monolithic model that must do everything needs enormous resources to run at all. Fifty specialised brains, each excellent at a narrow domain, can be loaded, unloaded, and combined according to the task and the hardware available. A small machine runs the brains it needs for the work in front of it and swaps them as the task changes. A large machine keeps more of them resident and reaches for them faster. The same fifty-brain design serves the laptop and the server, which is precisely what a doctrine meant for the real, uneven world requires.

I am deliberate about not overstating any single machine's reach. We detect the hardware at runtime, seal a profile of what is present, and scale accordingly, and we never claim a feature runs on silicon that cannot run it. The larger open models, up to the seventy-billion-parameter class, run on the workstation through hybrid processor and graphics-card offload, not served live from a single mid-range card, and I say so plainly rather than implying otherwise. That honesty is not a caveat bolted onto the marketing. It is part of the engineering, because a resilient system that lies to its operator about its own limits is just one surprise away from the failure it claimed to have designed out.



The Mickai pantheon.

Provenance as Evidence You Can Defend

The strongest economic argument for offline-first is not the meter you avoid. It is the evidence you accumulate. A system that seals every consequential action locally, under a recognised post-quantum standard, builds a defensible record of what it did and why, continuously, whether or not it was connected. For any operator who may one day have to answer to a regulator, a court, an auditor, or a board, that record is not overhead. It is an asset, and it is one that rented systems struggle to provide on the operator's own terms.

Think about what an audit normally costs when the trail lives in someone else's cloud. You request access, you wait, you receive what they choose to give you, and you trust that it was not altered. Each of those steps is a dependency and a vulnerability. When the audit record is sealed at the edge with a standard signature and anchored to an independent chain, the burden of proof moves in your favour. You are not asking a provider to vouch for you. You are holding cryptographic evidence that stands on its own, made at the moment the action occurred and verifiable by anyone with the public standard, not just the vendor who wrote it.

The cheapest insurance against the day you must prove what your system did is to have it prove every action as it happens.

This is also where the intellectual property behind the system becomes relevant to the argument rather than a footnote to it. The approach is protected by 101 filed UK patent applications, covering roughly 2,234 claims, owned by Mickai LTD, with the named inventor being Mickarle Wagstaff-Irons. I say filed rather than granted because that is the accurate status, and accuracy is the whole point of a

book about provenance. The portfolio matters here because it signals that the mechanisms described are specific, documented, and defensible, not vague gestures at resilience.

Put the pieces together and the economic case sharpens. Owning the hardware removes the meter. Specialised local brains keep the work going through outages that would stop a rented system cold. Sealing and anchoring turn every disconnected action into defensible evidence. Each of these has a cost, and each returns more than it costs to an operator whose work is heavy, sensitive, or exposed to the contested edge. The doctrine is not charity to the disconnected. It is a sound investment for anyone who cannot afford to be switched off.

The numbers will differ for every operator, and I would distrust any single figure that claimed to settle it. What does not differ is the shape of the argument. Up-front cost, durable ownership, near-zero marginal thinking, continuous defensible evidence, and survival through the moments that matter most. Weigh that against perpetual rent and a single fragile link, and for a serious operator the balance tilts hard toward owning the mind in the room.

PART FOUR: PUTTING THE DOCTRINE TO WORK

This part turns the argument into concrete instructions for builders, buyers, and the people who write the rules.

For the People Who Build Systems

If you build intelligent systems, the doctrine asks one discipline of you above all others. Make local the default and the network the enhancement, from the first line of the design, not as a retrofit. The order matters because the path of least resistance always pulls toward the cloud, and unless local-first is the foundation you stand on, it quietly erodes into a neglected fallback. Decide it early and defend it, because every shortcut later will be a shortcut back to dependency.

Audit your dependencies honestly. For every feature, ask what it truly requires, and refuse to let it acquire a connection to the network that it does not actually need. You will be surprised how many capabilities were leaning on the link out of habit rather than necessity. Cut those wires deliberately so that each function fails alone, and you will have built graceful degradation into the structure rather than hoping it emerges. The work is tedious. It is also where resilience actually lives.

Every dependency you add is a question you are choosing to answer with someone else's uptime.

Favour specialisation over a single enormous model. A focused brain that does one thing excellently can run on hardware a person actually owns, and a system of such brains scales across a hardware range in a way a monolith never will. This is the choice that makes local-first practical rather than aspirational. Fine-tune and specialise capable open foundations into narrow, excellent components, and you trade the impossible dream of one model that runs everywhere for the achievable reality of many that each run somewhere real.

Build accountability in at the edge. Seal consequential actions locally, at the moment they occur, using recognised standards rather than schemes of your own invention, and design a path to anchor that evidence independently when connectivity returns. Do not wait for the network to make your system trustworthy. A system that can prove what it did while disconnected is a system you can deploy into the hard places, and the hard places are where intelligence is worth the most.

Finally, be honest about limits in the software itself. When a capability needs more hardware or a connection it does not have, say so plainly and tell the operator what would change it. A system that degrades loudly and clearly is more trustworthy than one that fails silently and looks broken. Honesty is not the opposite of confidence. It is what confidence looks like when it has nothing to hide, and it is

the difference between a tool people rely on and one they learn to distrust.



The Mickai pantheon.

For the People Who Buy and Deploy Them

If you buy or deploy intelligent systems, you hold more power than you use. The market builds what buyers reward, and for years buyers have rewarded the demo on the strong link and never asked the harder question. Change what you ask for and you change what gets built. So make the hard question the first question, before the features, before the price. What does this system do when the link goes down? If the answer is evasive, you have learned the most important thing about it.

Insist on seeing the degraded case, not just hearing about it. Ask the vendor to pull the network in front of you and show you what survives. A system built offline-first will demonstrate this gladly, because the degraded case is its proud case, not its embarrassing one. A system that hesitates, or asks to reschedule, or shows you a thin cache and an apology, has told you where its intelligence really lives. The demonstration is cheap to ask for and impossible to fake.

Never buy an intelligent system on the strength of how it behaves when everything is going right.

Weigh the full cost over the life of the system, not the attractive figure at the start. Per-query pricing flatters the early months and punishes the later years, especially for heavy and sustained use. Owning the hardware costs more up front and far less over time, and it removes a meter that someone else controls. Do the arithmetic across the whole expected lifetime, and include the price of the outages you are pretending will not happen, because in your real environment they will.

Ask where your data and your audit trail live, and who can reach them. If the answer is a cloud you do not control, then your context, your memory, and your evidence are all held under terms you did not write and can be changed without your say. For sensitive work, that is a sovereignty problem you should not accept by default. Prefer systems that keep the context in the room, hold the memory locally, and seal their evidence where you can stand behind it. Ownership of the mind should come with ownership of what it knows.

And be wary of the word sovereign when it is used loosely, including by people like me. Press for specifics. Does the intelligence actually run on hardware you own? Does it actually work disconnected? Is the audit actually sealed locally under a recognised standard? Sovereignty is a property you can test, not a slogan you can buy, and a vendor who welcomes the testing is worth more than one who only offers the word.

For the People Who Write the Rules

If you write policy, procurement standards, or regulation, you shape the floor that everyone else builds on. Right now that floor quietly assumes connectivity, and so it permits, and even rewards, systems that vanish under stress. You can change that by making resilience under degraded conditions a stated requirement rather than a hoped-for virtue. Ask for the contested case in the specification, and the market will start building for it, because the market builds what the rules demand.

Treat the ability to function disconnected as a property worth certifying. For anything that operates in critical infrastructure, in healthcare, in defence, or anywhere a sustained outage costs lives or livelihoods, the question of what survives the loss of the link should be a formal part of how the system is assessed. A certificate that a system keeps its core functions through a contested period is worth more than a page of promises about uptime, because the contested period is precisely when promises tend to break.

Rules that only test systems on the easy days will keep producing systems that fail on the hard ones.

Value auditability that does not depend on a vendor's continued cooperation. When an audit trail is sealed at the edge under a recognised standard and anchored independently, an investigator can verify what happened without asking permission from the company under investigation. That property is a gift to anyone whose job is oversight, and policy can encourage it by preferring systems that produce their own defensible evidence over those that ask you to trust a remote log. Recognised standards such as the NIST post-quantum signature schemes give you a stable basis to point to.

Reward genuine sovereignty and refuse to be impressed by the word alone. Write tests into your procurement that distinguish a system that truly runs and audits on owned hardware from one that merely describes itself in sovereign language while depending on a distant cloud. The distinction is testable, and policy that tests it will steer public money toward resilience and away from fragility dressed up in fashionable words. The cost of getting this wrong is paid in the outages that hit hardest where help is thinnest.

I will close where I began. The link will go down. That is not pessimism, it is planning, and the difference between the two is whether you built for it. An intelligence that keeps thinking through the outage, that degrades gracefully rather than collapsing, that seals its actions where you can prove them, and that runs on hardware you own, is not a luxury for a contested world. It is the baseline that world deserves. Build to that baseline, buy to it, and write it into the rules, and intelligence stops being a privilege of good connectivity and becomes something people can actually rely on when it counts.



The Mickai pantheon.

APPENDIX · ABOUT THE AUTHOR

Micky Irons

Founder and chief executive of Mickai LTD (Companies House 17166618, registered office 20 Wenlock Road, London, N1 7GU) and named inventor on the Mickai SIOS patent corpus: 101 filed UK patent applications, around 2,234 claims. Trade mark Mickai registered at UK00004373277.

Profiles

mickai.co.uk

crunchbase.com/person/micky-irons

linkedin.com/in/mickyirons

© 2026 Mickai LTD. Set in Inter Tight and Inter Black. Brand voice audited; zero violations at publish.

References and further reading

- Kleppmann, M. et al. Local-First Software: You Own Your Data, in spite of the Cloud. Onward! ACM, 2019.
- National Institute of Standards and Technology. FIPS 204: Module-Lattice-Based Digital Signature Standard (ML-DSA). NIST, 2024.
- Nakamoto, S. Bitcoin: A Peer-to-Peer Electronic Cash System. 2008.
- Nygard, M. Release It! Design and Deploy Production-Ready Software, 2nd edition. Pragmatic Bookshelf, 2018.
- Hennessy, J. and Patterson, D. Computer Architecture: A Quantitative Approach, 6th edition. Morgan Kaufmann, 2017.
- Bernstein, D. and Lange, T. Post-Quantum Cryptography. Nature, vol. 549, 2017.
- Tanenbaum, A. and van Steen, M. Distributed Systems: Principles and Paradigms, 3rd edition. 2017.
- Schneier, B. Click Here to Kill Everybody: Security and Survival in a Hyper-connected World. Norton, 2018.