

# AI & Nearshore Teams: Who Gets Replaced and Why

A research-backed model of AI in nearshore software teams: who gets replaced, how incentives change, and what this means for CTOs building reliable pipelines.

TeamStation AI Research Paper Draft v1.0  
November 2025

## Authors

Lonnie McRorey\* — CEO & Co-Founder, TeamStation AI  
TeamStation AI Research Division

\*Corresponding author: [lonnie@teamstation.io](mailto:lonnie@teamstation.io)

## Table of Contents

### Abstract

### 1. Introduction

- 1.1 From job loss debates to pipeline design
- 1.2 What we study

### 2. Model

- 2.1 Agents, information, and effort choice
- 2.2 Project success technology
- 2.3 Contracts, timing, and the principal's problem
- 2.4 The role of AI
- 2.5 Expected cost under policy  $x$
- 2.6 The incentive derivative

### 3. Incentive Structure and Optimal Placement of AI in Sequential Teams

- 3.1 Why the last position is the most replaceable
- 3.2 Why the middle is structurally protected
- 3.3 Why the first position sits between the extremes
- 3.4 Why the optimal policy is probabilistic

### 4. Replacement Probabilities and Wage Effects under the Optimal Policy

- 4.1 Replacement probabilities decline as you move backward through the team

- 4.2 Wage effects rise at the start and in the middle
- 4.3 The optimal policy raises wages while lowering cost
- 4.4 Wage compression is a feature, not a bug
- 4.5 Replacement probabilities are interior, not corner solutions

## **5. Managerial Implications for Nearshore Engineering Pipelines**

### **6. Extensions and Directions for Future Research**

#### **Conclusion**

#### **Glossary of Terms**

#### **Mathematical Definitions and Notation**

#### **Data, Methodology, and Provenance Statement**

#### **Legal and Intellectual Property Notice**

## **Abstract**

Teams do not operate as isolated job titles. They work as linked stages where each effort choice changes what the next person believes is possible. Once AI shows up inside that chain, the question is no longer a blunt prediction about job loss. The real question is which positions get substituted, how often that substitution occurs, and what those choices do to the incentives of everyone who remains.

We build a sequential team production model in which each human worker selects effort ( $e_i$  equals 1) or shirking ( $e_i$  equals 0), pays a personal cost  $c$  greater than zero for effort, and sees only the effort choice of the person immediately before them. Project success depends on the count of workers who exert effort. If  $k$  workers try, the project succeeds with probability  $p_k$ , and these probabilities satisfy strict complementarity. For all  $k$  less than or equal to  $n$  minus 2:

$$p_{k+2} - p_{k+1} \text{ is greater than } p_{k+1} - p_k.$$

In simple terms, each new unit of effort adds more value when most of the chain is already engaged.

AI agents enter the model as reliable substitutes. They always exert effort, cost the principal  $c$ , and create no moral hazard. The principal, interpreted as a CTO or a system like Axiom Cortex, chooses success contingent wages  $w_i$  for humans and a probabilistic AI placement policy  $x$ . Each  $x_i$  lies between 0 and 1, and total AI capacity must satisfy the constraint that the sum of  $x_i$  is less than or equal to 1.

Given any  $x$ , the least cost way to sustain full effort from all humans yields

$$w_i^x \text{ equals } c \text{ divided by } (p_n - \zeta_i^x)$$

where  $\zeta_i^x$  is the probability the project still succeeds when worker  $i$  shirks. This term depends on AI placements downstream. A higher value makes shirking safer and pushes wages up.

Three forces determine whether raising  $x_i$  reduces or increases the principal's total expected cost. The first is direct cost saving. Replacing a human avoids the expected payment  $p_n w_i^x$  and pays the constant cost  $c$ . The second is direct incentive cost. Increasing  $x_i$  raises  $\zeta_i^x$ , which increases  $w_i^x$ . The third is indirect incentive cost. Changing  $x_i$  shifts the shirking temptation  $\zeta_k^x$  for all workers upstream, which raises their wages as well.

These interactions create an uneven structure inside the pipeline. The middle position becomes essential because removing it breaks the information link that peer monitoring relies on. The end position is most exposed to AI because  $\zeta_n^x$  equals  $p_{\{n \text{ minus } 1\}}$  and remains unchanged when AI is added, which lets the principal capture the largest savings without raising incentives at that stage. The first position sits between these two. Replacing the first saves less and disturbs fewer incentives.

The optimal AI policy is probabilistic. No position is always automated, and no position is always protected. The principal sometimes leaves part of the AI capacity unused because the risk of full human failure helps keep wages disciplined.

Under the optimal policy, wages for the first and middle positions rise, the wage for the end stays constant, and the internal wage spread narrows. Put simply, the correct use of AI lifts the bottom of the wage structure, leaves the top in place, and reduces internal inequality.

**We link these results to nearshore engineering pipelines** and the TeamStation AI platform. Axiom Cortex measures cognitive discipline and task fit, while Nebula orchestrates which steps fall to humans or AI across distributed projects. The central lesson is that AI placement is not only a technical decision. It is an incentive design decision that determines whether teams remain reliable once automation becomes part of the workflow.

## 1. Introduction

The conversation around AI and work often falls into a strange loop. People ask whether machines will replace developers, analysts, testers, or designers, as if the labor market were a collection of disconnected seats waiting to be swapped out. Actual teams do not function that way. A team is a chain of dependencies. What happens at one step shapes the beliefs, risks, and expectations at the next. The structure of those dependencies is what determines whether AI improves output or quietly breaks the system.

Nearshore engineering teams make this even more visible. The work moves across time zones, talent pools, and communication layers, and the reliability of each stage depends on what came before. Anyone who has built delivery pipelines in practice has seen the pattern. A strong start creates momentum. A weak start creates hesitation. And once hesitation enters, every downstream actor starts asking a simple question: does my effort actually matter if the earlier steps have already drifted off course.

That question lives at the center of this paper. We study a team that works in sequence and where each human chooses effort or shirking. Each worker pays a personal cost  $c$  when they contribute real effort. Each one sees only the effort choice of the person immediately before them. This limited visibility creates a familiar form of peer monitoring. If a worker sees effort, they feel safe enough to exert effort themselves. If they see shirking, they rationally reduce their own effort because the project already looks compromised.

The principal, interpreted as a CTO or an operational platform like Axiom Cortex, designs a contract that uses success contingent wages to keep every worker engaged. The project succeeds with probability  $p_k$  when exactly  $k$  workers exert effort. These probabilities satisfy strict complementarity. For all  $k$  less than or equal to  $n$  minus two, the condition

$$p_{k+2} - p_{k+1} > p_{k+1} - p_k$$

captures the idea that effort is more valuable when most of the chain is already contributing. This is something every engineering leader has felt even without the math. A late stage contribution from a strong engineer means little if the earlier stages were already drifting. The same contribution has enormous leverage if the earlier stages were disciplined and aligned.

AI enters this structure as a reliable unit of effort. It always exerts effort, costs  $c$ , and introduces no moral hazard. It also breaks none of the sequential logic. The principal chooses a probabilistic AI placement plan  $x = (x_1, \dots, x_n)$  with the capacity constraint that the sum of  $x_i$  is at most one. Each  $x_i$  is the probability that the worker in that position is replaced by AI. The goal is simple in theory: reduce expected cost while keeping every human in the chain willing to exert effort.

What makes the problem interesting is that AI does not simply replace a cost. It changes the entire incentive surface. AI downstream from a worker makes shirking safer for that worker because the project may still succeed even if they slack. That value is captured in  $\zeta_i^x$ , the success probability when worker  $i$  shirks under policy  $x$ . Wages that keep workers honest take the form

$$w_i^x \text{ equals } c \text{ divided by } (p_n - \zeta_i^x).$$

This expression shows why AI placement becomes a subtle design problem. Increasing AI in one position can lower cost directly but raise wages elsewhere. Some positions amplify this effect, and some dampen it.

This paper studies how those forces combine, where AI belongs inside a sequential workflow, and how a cost minimizing principal balances direct savings against incentive pressure. The concepts map naturally onto nearshore engineering because the same pattern appears in distributed delivery. One stage shapes another. Beliefs about reliability shift decisions. A single break in the chain forces everyone downstream to discount their own effort.

What follows is a model that captures these mechanics. It is abstract on purpose. The simpler the structure, the easier it becomes to see the underlying incentives clearly. And once those incentives are visible, the decisions about AI placement in real engineering pipelines become easier to reason about and

easier to systematize. That is the goal of this work. Not predictions about job loss. A design surface for where AI fits inside real teams and why certain seats are replaceable while others carry structural weight.

## 1.1 From job loss debates to pipeline design

Most conversations about AI fall into a habit that feels tidy but hides the real structure underneath. People ask whether the technology will take a job, or a set of jobs, or eventually sweep through entire categories. It sounds reasonable until you look at how teams actually work. Real teams are sequences. One step hands work to the next, sometimes cleanly, sometimes with a bit of grit that the next person must smooth out. The point is that the value of any individual step depends on what came before it. Job loss arguments miss that completely.

The shift from a job lens to a pipeline lens changes the stakes. Once you look at a team as a set of contingent signals, the question snaps into focus. You are no longer asking who gets replaced. You are asking where inside the workflow an AI unit of effort helps and where it quietly distorts the incentives that hold the rest of the chain together. That distortion is subtle. It creeps in through belief changes rather than through technical capability, and it is easily overlooked if you treat the pipeline as a flat list of roles.

**Anyone who has worked in nearshore delivery** has seen this in the wild. A strong opening developer gives the next engineer a clear structure to build on. A shaky start does the opposite. It creates hesitation, and hesitation spreads. People downstream start to wonder whether their own effort will matter. And once that question enters the room, effort drops. It is not laziness. It is a rational adjustment to a weaker probability of success. That adjustment is the reason peer monitoring works in the first place.

AI slips into this picture with a very particular kind of weight. It is reliable, which is useful, and it is predictable, which is dangerous. Once AI occupies a downstream slot, the worker just above it realizes that shirking is less costly than before. The system might still succeed. That small change in the conditional probability of success is exactly what we capture with  $\zeta_n^x$  in the model. It is the success probability given shirking. And once that number rises, the wage that keeps the worker honest rises too. That is the core mechanism. It is not philosophical. It is mechanical.

There is a more intuitive way to see it. Think of a relay race where one runner never tires and never drops the baton. If you run just before that person, you know the anchor will salvage almost anything you hand off. This gives you a small but meaningful temptation to ease off. A manager who wants everyone to run at full speed suddenly has a new incentive problem, even though the anchor is perfect. In distributed engineering teams, this effect appears every day. You can see it in how code reviews change when the reviewer is strict versus when the reviewer is permissive. Downstream reliability shifts upstream behavior.

This is why the old job loss framing fails. It treats all positions as interchangeable. But the incentive impact of AI depends heavily on placement. Some seats cannot absorb the change in belief without breaking the chain. Some seats respond mildly. And one seat, the end of the pipeline, barely changes at all because its  $\zeta_n^x$  already equals  $p_{\{n \text{ minus } 1\}}$ . The model formalizes that intuition, but the intuition itself comes from watching how teams behave under pressure.

This paper steps away from abstract predictions and moves toward a design surface for AI placement. The question is not who gets replaced. The question is what happens to the incentive gradient when you insert a perfectly reliable worker at a specific point in the sequence. Once you see the problem that way, the rest of the analysis follows naturally.

## 1.2 What we study

This paper examines a principal who manages a team arranged in sequence, where each worker hands the project to the next. The team has  $n$  human workers. Each worker  $i$  chooses an effort level  $e_i$  that is either 1 for effort or 0 for shirking. Putting in effort costs  $c$  greater than zero. A worker sees only the effort choice of the worker immediately before them, which creates a narrow but realistic information channel similar to what nearshore pipelines face when work moves between stages or time zones.

The project succeeds with probability  $p_k$  when exactly  $k$  workers exert effort. These probabilities rise with  $k$  and also satisfy strict complementarity. For all  $k$  less than or equal to  $n$  minus two:

$$p_{k+2} - p_{k+1} \text{ is greater than } p_{k+1} - p_k.$$

This condition captures a feature teams know well. When most of the chain is working hard, the next unit of effort is unusually valuable. When the chain is weak, the same unit of effort makes a smaller difference. Anyone who has spent nights salvaging failing delivery pipelines will recognize the pattern.

The principal selects two objects. First is a vector of success contingent wages  $w = (w_1, \dots, w_n)$  that make humans willing to exert effort. Second is an AI replacement policy  $x = (x_1, \dots, x_n)$  where each  $x_i$  lies between 0 and 1 and the total capacity satisfies the constraint that the sum of  $x_i$  is at most 1. The value  $x_i$  is the probability that worker  $i$  is replaced by AI. AI units behave as perfectly reliable workers. They always exert effort, create no moral hazard, and cost  $c$ . The principal may use them to reduce expected cost if the wage pressure on humans becomes too high.

We characterize the equilibrium in which all humans choose effort. The object that determines their incentive is  $\zeta_i^x$ , the success probability when worker  $i$  shirks under replacement policy  $x$ . The  $\zeta$  term rises when AI is added downstream because AI can rescue the project even when humans fail. The least cost wage that keeps worker  $i$  willing to exert effort is

$$w_i^x \text{ equals } c \text{ divided by } (p_n - \zeta_i^x).$$

This expression is central to the analysis. It shows that wages and AI placement are tied through incentives, not just through substitution cost. A change in  $x_i$  may help the principal directly but cause an increase in wage pressure elsewhere. The net effect can be positive or negative depending on the structure of the sequence.

The rest of the paper studies four questions. First, where should AI be placed to minimize expected cost while preserving full effort from all remaining humans. Second, which positions face the highest replacement probability under the optimal plan. Third, how that plan changes the wage distribution inside

the team. Fourth, whether the principal sometimes prefers a stochastic replacement rule instead of a deterministic one.

These questions matter for engineering teams because the model isolates a mechanism leaders see but rarely formalize. Incentives shift when reliable automation appears downstream. Some positions harden. Others become fragile. Understanding this structure helps CTOs place AI inside real nearshore pipelines without weakening the behavioral foundation that keeps those pipelines reliable.

## 2. Model

The model describes a team arranged in sequence. Each position handles a stage of work and then hands it to the next one. This matches how nearshore delivery actually behaves. Time zone alignment helps, but the incentive mechanics do not change with geography. Effort in one place shapes expectations everywhere else.

The team has  $n$  human workers, indexed from 1 to  $n$ . Worker  $i$  chooses effort  $e_i$  that is either 1 or 0. Effort costs  $c$  greater than zero. Shirk costs nothing. Each worker sees only the effort choice of the previous one. That narrow visibility is the entire backbone of peer monitoring. It also introduces fragility, since a single observed shirk can ripple through the rest of the line.

The principal, interpreted as a CTO or as a system assigning tasks inside a nearshore pipeline, wants every human to choose effort. The principal uses two levers to make that happen. The first is a vector of success contingent wages  $w = (w_1, \dots, w_n)$ . The second is an AI placement policy  $x = (x_1, \dots, x_n)$ , where each  $x_i$  lies between zero and one and the total use of AI satisfies the constraint that the sum of  $x_i$  is at most one. Each  $x_i$  is the probability that worker  $i$  is replaced by AI.

### 2.1 Agents, information, and effort choice

A human in position  $i$  sees only  $e_{\{i \text{ minus } 1\}}$ . This limited information structure is not an abstraction. It mirrors real handoffs. Engineers and analysts do not see everything upstream. They see the output of the person right before them. And they infer reliability from that single signal.

Workers behave as follows. If they choose effort, they pay  $c$ . If they shirk, they pay zero. They compare the expected payoff from exerting effort to the expected payoff from shirking. That comparison depends on the probability the project succeeds in each case. When a worker shirks, the project may still succeed because others downstream exert effort, or because AI takes over one of the remaining slots. This leads to the key object of the model:  $\zeta_i^x$ , the probability of success when worker  $i$  shirks under policy  $x$ .

This single value encodes the entire incentive distortion created by AI downstream. When AI appears after  $i$ ,  $\zeta_i^x$  rises. Shirk becomes safer. Exerting effort becomes relatively less valuable. And the wage needed to keep  $i$  honest increases.

## 2.2 Project success technology

Project success depends on the number of workers who exert effort. If exactly  $k$  workers exert effort, the success probability is  $p_k$ . These probabilities satisfy strict complementarity. For all  $k$  less than or equal to  $n$  minus two:

$$p_{k+2} - p_{k+1} \text{ is greater than } p_{k+1} - p_k.$$

Teams experience this all the time. A good sequence compounds. A bad sequence drags. Effort delivered into a strong pipeline has disproportionate leverage. Effort delivered into a weak pipeline often evaporates on contact.

These  $p_k$  values are exogenous. They describe technology. They do not depend on wages or AI placement. They do, however, shape how incentives propagate through the team.

## 2.3 Contracts, timing, and the principal's problem

The principal commits to wages  $w_i$  and to the AI placement probabilities  $x_i$ . Wages are paid only when the project succeeds. That means the payoff to worker  $i$  is either  $w_i$  or zero depending on project completion. Worker  $i$  compares two expected payoffs:

1. Effort:  $p_n w_i - c$
2. Shirk:  $\zeta_i^x w_i$

The principal wants  $p_n w_i - c$  to be at least  $\zeta_i^x w_i$ . Solving that inequality yields the least cost wage that keeps  $i$  willing to exert effort:

$$w_i^x \text{ equals } c \text{ divided by } (p_n - \zeta_i^x).$$

This expression contains all of the economic action. The term  $(p_n - \zeta_i^x)$  is the difference in success probability between effort and shirking. It is the incentive margin. Once AI begins to occupy downstream seats, that incentive margin shrinks. And wages must rise to compensate.

This is how AI changes team economics. It does not change the project technology  $p_k$ . It changes the incentive surface that sits beneath the technology.

## 2.4 The role of AI

AI in this model is simple. It always exerts effort. It costs  $c$ . It generates no moral hazard. And it is treated as a unit that fills a position with probability  $x_i$ .

This is deliberate. For nearshore teams, AI does not replace human cognition. It replaces the risk that a human shirks or delivers inconsistent work. AI is a stabilizer. That stabilizer is helpful, but only when placed in the right slot. Put it in the wrong slot and it removes the strategic uncertainty that keeps upstream humans disciplined.

That is why the principal does not always spend the full AI budget. Sometimes the threat of full human failure is more valuable than an additional reliable step.

## 2.5 Expected cost under policy x

The principal's expected cost includes two parts. When a human occupies position  $i$ , the expected wage payment is  $p_n w_i^x$ . When AI occupies it, the cost is  $c$ . Because  $x_i$  is the probability of AI replacement, the expected cost contribution of position  $i$  is:

$$x_i c + (1 - x_i) p_n w_i^x.$$

Summing across all positions yields the principal's total expected cost under policy  $x$ .

This cost includes direct savings from replacing high wage humans with AI. It also includes indirect costs from the incentive distortions caused by that same replacement. The model identifies the balance point where the two forces equalize.

## 2.6 The incentive derivative

To understand how the principal chooses  $x_i$ , we differentiate the expected cost with respect to  $x_i$ . The derivative decomposes into three pieces:

1. Direct cost saving:  $(p_n w_i^x - c)$
2. Direct incentive cost:  $(1 - x_i) p_n \times (\text{partial derivative of } w_i^x \text{ with respect to } x_i)$
3. Indirect incentive cost:  $\sum_{k < i} (1 - x_k) p_n \times (\text{partial derivative of } w_k^x \text{ with respect to } x_i)$

Each term has a clear meaning. Direct savings push the principal to increase  $x_i$ . Incentive costs push the principal to reduce  $x_i$ . The sign of the derivative tells the principal whether increasing AI exposure at position  $i$  lowers or raises expected cost.

These expressions are structural. They hold regardless of geography, wage rate, or team seniority. This is why the model applies cleanly to nearshore engineering. The human logic is identical. The only difference is the practical value of time zone alignment, which affects project coordination but not the incentives inside this model.

# 3 Incentive Structure and Optimal Placement of AI in Sequential Teams

Teams arranged in sequence do not respond symmetrically when automation enters the line. The effect of replacing one position depends entirely on how beliefs and incentives propagate upstream. This section characterizes those propagation effects and shows which positions are structurally exposed to AI, which are protected by the incentive architecture, and why the principal often chooses a probabilistic replacement policy rather than a deterministic one.

The expressions developed in earlier sections give us the core object:

$$w_i^x \text{ equals } c \text{ divided by } (p_n - \zeta_i^x).$$

This links every worker's wage to the success probability when they shirk. Once AI appears downstream,  $\zeta_i^x$  increases, and the denominator shrinks. Incentives tighten. Wages must rise. The more downstream the AI, the more pronounced the ripple upstream.

This dynamic is identical in nearshore engineering pipelines. A single guaranteed-reliable stage downstream allows upstream engineers to interpret failure as less damaging, which raises the temptation to hold back. Leaders see this in code reviews, QA triage, and staging cycles. When a step becomes "too reliable," upstream discipline weakens unless incentives counterbalance the shift.

We now examine the derivative of the principal's expected cost with respect to  $x_i$ . The sign of that derivative determines whether increasing AI at position  $i$  is desirable.

## 3.1 Why the last position is the most replaceable

The end of the pipeline behaves differently from every other point in the sequence. When worker  $n$  shirks, the project succeeds with probability  $p_{\{n-1\}}$ . Adding AI after them is impossible because there is no "after." This means  $\zeta_n^x$  equals  $p_{\{n-1\}}$  regardless of  $x$ . No AI placement anywhere in the chain alters their shirking margin.

From the wage equation, this implies:

$$w_n^x \text{ equals } c \text{ divided by } (p_n - p_{\{n-1\}}),$$

which is constant across all policies.

Because wages for the end never rise in response to AI elsewhere, replacing them yields the cleanest savings. The principal avoids paying  $p_n w_n^x$  and instead pays  $c$ . Since  $w_n^x$  does not depend on  $x$ , this is all gain and no incentive distortion. That is why, in both the mathematical model and real distributed pipelines, the final stage tends to be the most suitable location for automated effort.

In nearshore engineering, this often corresponds to QA validation, data aggregation, error-checking passes, or final documentation transforms. These steps are structurally tolerant to automation because no worker depends on observing them before making their own effort decision.

## 3.2 Why the middle is structurally protected

Replacing a middle position disrupts the informational link that peer monitoring depends on. Worker  $i$  observes  $e_{i-1}$ . Worker  $i+1$  observes  $e_i$ . If  $i$  becomes AI, both neighbors experience a different incentive landscape. Upstream workers suddenly face a higher  $\zeta_k x$  because the AI downstream guarantees some baseline probability of success even when they shirk. Downstream workers lose a human signal they relied on to update expectations.

Mathematically, the zeta function for  $k < i$  rises sharply when AI sits at  $i$  because:

- the AI at  $i$  always exerts effort
- the downstream chain becomes more reliable
- the conditional success probability given shirking increases

This causes the partial derivative of  $w_k x$  with respect to  $x_i$  to spike for all  $k < i$ .

The middle is where these effects overlap the most. It affects both directions. Replacing this seat raises wage pressure upstream and weakens monitoring downstream. This is why the optimal solution avoids placing AI deterministically in the middle. In both the formal model and real engineering teams, the center of the workflow carries structural weight that cannot be automated without penalty.

## 3.3 Why the first position sits between the extremes

The first worker does not observe anyone. They carry no peer monitoring load. Replacing them with AI avoids paying the expected wage  $p_n w_1 x$  and introduces no downstream informational loss. However, replacing the first raises  $\zeta_2 x$ , which increases  $w_2 x$ . This cost is milder compared to the cascading distortions caused by replacing a middle worker.

The end worker is the easiest to replace.

The first worker is moderately substitutable.

The center is the least substitutable.

This ordering emerges directly from the structure of the zeta function and the complementarity condition  $p_{k+2} - p_{k+1} > p_{k+1} - p_k$ . The incentive margin is most sensitive in the middle because effort there has the largest effect on project success when the chain is functioning well.

Nearshore engineering reflects this cleanly. Strong leads, mid-pipeline integrators, and architecture maintainers hold the system together. Replace them incorrectly and teams begin to drift. Replace the first or last and the drift does not occur in the same way.

### 3.4 Why the optimal policy is probabilistic

Even when a position is attractive for AI, the principal may choose not to replace it deterministically. The reason lies in the form of the incentive derivative:

1. direct cost saving
2. direct incentive cost
3. indirect incentive cost

The total effect can be convex or concave depending on how zeta shifts. When the relationship is non-linear, the cost minimizing solution is a probability between zero and one. **Using AI stochastically** preserves enough uncertainty to maintain upstream discipline while still capturing cost savings in expectation.

This is one of the counterintuitive results of the original model and it helps CTOs reason about workflow design. The correct use of automation in team sequences is rarely “all or nothing.” It is often an exposure level that preserves the incentive gradient without flattening it.

In practice, nearshore organizations do this implicitly when they keep some stages strictly human even though they possess AI tools capable of filling them. They maintain a living incentive structure rather than a sterile technical pipeline.

## 4 Replacement Probabilities and Wage Effects under the Optimal Policy

Once the principal balances direct cost savings against the incentive distortions created by AI, the structure of the optimal policy becomes visible. Some positions are replaced with high probability, some with low probability, and none are replaced with certainty. The resulting wage pattern is not an accident. It is a mechanical consequence of how  $\zeta_i^x$  propagates through the team.

This section characterizes those probabilities and the wage effects they generate.

### 4.1 Replacement probabilities decline as you move backward through the team

The end worker remains the most exposed to AI. Their shirking margin does not depend on the placement of AI elsewhere. Formally,

$$\zeta_{n^x} \text{ equals } p_{\{n \text{ minus } 1\}}.$$

This fixes their wage at

$$w_{n^x} \text{ equals } c \text{ divided by } (p_n \text{ minus } p_{\{n \text{ minus } 1\}}).$$

Because this expression remains constant across all  $x$ , replacing the last position saves  $p_n w_{n^x}$  minus  $c$  with no incentive cost anywhere else. This yields a replacement probability for the final position that is largest among all positions in the sequence.

Moving backward, the effect weakens. Any replacement at position  $n$  minus 1 raises  $\zeta_{\{n \text{ minus } 1\}^x}$  and  $\zeta_{\{n \text{ minus } 2\}^x}$ , which raises  $w_{\{n \text{ minus } 1\}^x}$  and  $w_{\{n \text{ minus } 2\}^x}$ . The indirect pressure intensifies the further from the end the principal tries to automate.

This creates a clear pattern:

- high probability of replacement at the end
- moderate at the start
- lowest at the center

The center's protection is not philosophical. It is an equilibrium property. AI in the middle tightens incentives on too many workers at once.

In nearshore engineering languages, this matches what CTOs see with architecture and integration roles. These are often the least automatable not because AI cannot perform the tasks, but because the human link they provide stabilizes the incentive gradient inside a distributed workflow.

## 4.2 Wage effects rise at the start and in the middle, and stay constant at the end

Given the wage equation

$$w_{i^x} \text{ equals } c \text{ divided by } (p_n \text{ minus } \zeta_{i^x}),$$

any position whose  $\zeta_{i^x}$  rises under the optimal policy experiences an increase in wages.

The end worker's wage remains fixed, since their  $\zeta$  value never changes.

The first and middle workers experience the largest increases.

The reasons differ.

### For the first worker:

AI anywhere downstream increases  $\zeta_1^x$ , because the first worker expects the project to be rescued

even when they shirk. That increases  $w_1^x$ . This is intuitive. When the downstream chain becomes more reliable, the first worker's marginal contribution shrinks. They demand compensation for that loss.

**For middle workers:**

AI placed downstream raises  $\zeta_i^x$  sharply, because the probability of successful rescue increases most when the pipeline is already strong. With strict complementarity in the  $p_k$  sequence, the center is where that effect is strongest. The wage increase reflects that amplification.

**For the end worker:**

Their wage stays constant because neither  $\zeta_n^x$  nor the difference  $p_n - \zeta_n^x$  ever moves.

The internal wage structure therefore compresses.

The bottom rises.

The top does not move.

This narrowing appears even in real nearshore pipelines. When automation stabilizes late stage steps, juniors at the front and mids in the center require more disciplined coordination and thus command higher expected compensation. The end roles, already least sensitive to pipeline drift, retain stable economics.

### 4.3 The optimal policy raises wages for humans but lowers total cost

This apparent contradiction is one of the model's most important results.

When the principal installs AI correctly:

- some wages rise, because incentives must be restored
- expected total cost falls, because AI saves more at the end than it costs in wage adjustments elsewhere

The wage increases are targeted. They fall on positions whose incentives become fragile when the pipeline becomes more reliable. The savings come from replacing the seats where reliability matters most and incentives matter least.

Formally, if  $x^*$  is the optimal policy, then for any feasible alternative  $x\tilde{}$ :

$\text{ExpectedCost}(x^*) \leq \text{ExpectedCost}(x\tilde{})$

even though for several positions  $i$ :

$w_i^{x^*} > w_i^{x\tilde{}}$ .

This is not a trick. It is how sequential incentive systems behave.

You save on cost where incentives are flat, and you compensate where incentives are steep.

For CTOs operating nearshore pipelines, this means automation raises the bar on some roles and lowers cost overall only when placed at structurally insensitive positions. The model formalizes what experienced delivery leaders already do intuitively.

## 4.4 Wage compression is a feature, not a bug

A striking result is that the internal wage difference  $w_n^x - w_1^x$  shrinks under the optimal policy.

The end wage remains fixed.

The first worker's wage rises.

The middle rises even more.

This compression emerges because  $\zeta_i^x$  rises more sharply at positions where the marginal value of effort is highest. Complementarity in the  $p_k$  sequence makes the middle particularly sensitive. Incentives need support, so wages rise. The end requires none, so its wage is stable.

Compression has implications for nearshore engineering:

- junior and mid roles gain relative importance
- late stage roles become more standardizable
- the difference between high leverage and low leverage roles narrows

This matches what happens when AI begins to clean up late stage tasks like error checking, logging normalization, and documentation transforms. Human leverage migrates upstream.

## 4.5 Replacement probabilities are interior, not corner solutions

One of the most important results is that for most positions:

**$x_i^*$  lies strictly between 0 and 1.**

This interior solution arises from the non-linearity of incentive distortions.

A deterministic rule dulls the incentive margin too sharply.

A probabilistic one preserves enough uncertainty to sustain discipline.

The model shows that stochastic automation is not a theoretical artifact. It is an efficient mechanism for maintaining reliability in sequential pipelines. Real organizations approximate this when they apply automation inconsistently, by design. They keep humans occasionally in the loop, not because the AI is worse, but because the incentives downstream become too weak if the automation is too predictable.

In **nearshore delivery terms**, this is why **hybrid approaches outperform pure automation**. Keeping some tasks human with non-zero probability maintains effort quality across the entire chain.

## 5 Managerial Implications for Nearshore Engineering Pipelines

The results of the model draw a clean map for how CTOs should deploy automation inside distributed nearshore engineering teams. The math does not speak in metaphors, but the implications do. The structure of incentives in a sequential team tells us which roles benefit from automation, which roles must remain human for the system to hold, and how wages and expectations shift once reliable automation becomes part of the production chain.

The first implication is that **end stage roles are the most automatable**, not because the work is the simplest, but because their incentives are flat. In the model, the end worker's zeta value equals  $p_{n-1}$  regardless of where AI is placed. Their wage is fixed at  $w_{n-1}x$  equals  $c / (p_n - p_{n-1})$ . Replacing them carries no incentive distortion, which makes these positions clean automation targets. In nearshore pipelines, this corresponds to tasks like synthetic QA runs, safety validations, documentation transforms, or batch data checks. These steps have low informational spillback into the chain.

The second implication is that **the center of the pipeline is structurally protected**. Middle roles sit on the steepest part of the incentive gradient. Effort delivered at the center has the highest marginal effect when the sequence is strong. Removing a middle human breaks the informational link on which peer monitoring relies. Mathematically, the partial derivative of  $w_kx$  with respect to  $x_i$  for  $k < i$  spikes when AI occupies a central seat. This means the entire upstream wage budget inflates, making deterministic replacement at the center economically unattractive. In nearshore engineering, these are architecture, integration, sequence control, and key mid-layer implementation roles.

The third implication is that **first position automation is attractive but not as clean as end automation**. The first worker does not observe anyone, so the incentive structure has no upstream component. Replacing them avoids  $p_n w_1 x$  and disturbs only the incentives of the second worker. This makes the first position conducive to AI support, but it is not frictionless. It raises  $\zeta_2 x$  and therefore raises  $w_2 x$ . This is modest compared to the distortions caused by automating the center, but it is real. In nearshore pipelines, first role automation corresponds to automated scaffolding, project initialization, or base code generation.

The fourth implication is that **optimal automation is probabilistic**. A CTO who deploys automation deterministically at a given position weakens the incentive margin for all upstream workers. The model shows that the cost function often has an interior minimum, where the replacement probability  $x_i^*$  lies between zero and one. This is not a quirk of the math. It reflects a practical truth in distributed teams. Humans maintain effort because the environment retains some uncertainty. If automation becomes

perfectly predictable at a certain stage, the upstream incentive gradient flattens and wage or oversight costs rise. A stochastic or hybrid policy preserves enough uncertainty to keep upstream behavior aligned.

The fifth implication is that **internal wage compression is both predictable and beneficial**. The model shows that the correct use of AI raises wages at the beginning and middle of the chain while leaving the end wage fixed. Formally, wages at  $i < n$  rise because  $\zeta_i x$  increases with downstream automation. Wages at  $n$  remain at  $c$  divided by  $(p_n - p_{n-1})$ . This narrows the internal spread. Compression signals that responsibility is shifting upstream. It also signals that mid-tier roles are becoming more critical as automation absorbs the late pipeline variance.

The final implication is that the **principal sometimes leaves unused automation capacity on the table**. In the model, the constraint that the sum of  $x_i$  is at most one binds only at the optimum if the incentive distortions are weak. If incentive distortions are strong, the principal deliberately leaves some  $x$  capacity unused to preserve the threat of human failure. In real nearshore engineering operations, this appears when teams intentionally keep humans in loop stages like architecture review, integration design, and acceptance testing, even when AI tools are available. It is not reluctance. It is incentive design.

For US CTOs building nearshore pipelines, this section yields a simple map.

- Automate the end.
- Support the first.
- Protect the center.
- Use hybrid policies.
- Expect wage compression.

And preserve enough uncertainty that upstream effort remains disciplined.

These patterns arise from math, not management taste. They provide a template for building **stoichiastic** and heterogeneous cognitive architectures that reflect the underlying economics of effort and belief inside a distributed team.

## 6 Extensions and Directions for Future Research

The model we have developed is deliberately simple. That simplicity is what makes the incentive mechanics visible. But real nearshore engineering pipelines contain richer structures than a single linear chain. Time zones introduce small lags. Code dependencies create local complementarities. AI tools vary in reliability, speed, and context sensitivity. And human effort is not a binary variable. Extending the model along these lines opens clear paths for future work.

The first extension is to allow **heterogeneous effort costs**. In the baseline model, every worker pays the same cost  $c$  for effort. In nearshore teams, effort costs vary by seniority, by cognitive load of the task, and by how much context switching the position requires. Letting  $c_i$  differ across positions would shift the

wage equation to  $w_i^x$  equals  $c_i$  divided by  $(p_n - \zeta_i^x)$ . This preserves the incentive structure but changes the cost minimizing solution. Positions with lower  $c_i$  become more attractive for AI because the wage floor is smaller, while positions with high  $c_i$  become more sensitive to  $\zeta$  shifts. This would give CTOs a clearer map of which seniority bands are structurally stable under automation.

A second extension is to replace the strict complementarity condition with **task specific production functions**. In practice, different pipelines exhibit different shapes. Some have nearly linear returns to additional effort. Others display sharp returns only when context is fully aligned. The model would allow this by replacing the sequence  $p_1, \dots, p_n$  with any increasing success function  $F(k)$  that preserves the property that  $F(k+1) - F(k)$  is non-decreasing. A team whose architecture stage has unusually high leverage would produce a spike in the incentive gradient at that position. This could formalize why architectural roles remain the least automatable.

A third extension is to incorporate **multi stage visibility**. In the current model, worker  $i$  sees only  $e_{\{i-1\}}$ . Real engineering teams often see partial signals across several upstream steps. A worker may not know exactly what happened three positions earlier, but they may sense quality or reliability indirectly. Extending the model so that worker  $i$  observes a noisy signal of multiple upstream choices would preserve the core incentive logic while softening the sharp discontinuities at each position. This would make incentive distortions from AI replacement more diffuse but still measurable.

A fourth extension is to allow **AI reliability below one**. In the baseline model, an AI unit always exerts effort. In practice, different models and different tools have different failure rates. A tool that is reliable with probability  $r$  less than one would enter the  $\zeta$  function in a graded manner. That would shift the optimal replacement pattern. Highly reliable AI would still occupy the end of the sequence. Less reliable AI might be excluded entirely or used only probabilistically. This would match what CTOs already practice when they mix deterministic tools with human oversight loops.

A fifth extension is to incorporate **dynamic incentives and learning**. Workers in long lived nearshore pipelines form beliefs over time. If AI repeatedly rescues late stage failures, upstream workers update their expectations of project salvage. This turns  $\zeta_i^x$  into a dynamic variable. Wages would then depend on the entire history of automation, not just current placement. Extending the model to a repeated setting with Bayesian updating would capture how incentive gradients drift over time and why teams sometimes need periodic resets to restore discipline.

A sixth extension is to integrate **Axiom Cortex style cognitive heterogeneity**. Workers do not supply identical effort. They supply effort with distinct cognitive profiles: pattern matching, analytic depth, error resistance, and contextual reasoning. AI tools also differ across these axes. Embedding a multidimensional effort space into the model would let the principal choose not only where to place AI, but which type of AI to place. This matches the practical design questions CTOs face when mixing deterministic static analysis tools with generative models and human reviewers.

A final extension is to consider **parallel pipelines**. Many engineering organizations run work in branches that reconverge. Parallelization weakens local incentives because individual workers believe that other branches may salvage failures. Formalizing this would allow the model to describe hybrid architectures

common in CI/CD environments. It would also show when AI belongs in a branch and when it should sit only at convergence points.

Each of these extensions preserves the core logic:

$\zeta_i^x$  shapes incentives.

$w_i^x$  equals  $c$  divided by  $(p_n - \zeta_i^x)$ .

AI changes  $\zeta$  downstream and therefore changes the wage structure upstream.

The simplicity of the original model makes these extensions easy to integrate. They offer a clean research program for understanding how AI and human incentives interact in the nearshore engineering systems that US CTOs operate every day.

## Conclusion

This paper develops a simple structure with sharp consequences for how AI reshapes incentives inside sequential engineering teams. The math shows that the position a worker occupies in the pipeline determines how sensitive their effort is to downstream reliability. Once AI is introduced, those sensitivities become economic facts. They dictate which roles are easy to automate, which remain structurally human, and how wages adjust when reliability shifts.

The core result is clear. End positions are the most replaceable because AI does not distort their incentives. Middle positions are the least replaceable because they sit on the steepest part of the incentive gradient. First positions fall in between. The optimal automation strategy is rarely deterministic. It often assigns an interior replacement probability, which preserves enough uncertainty to keep upstream workers aligned. This generates a predictable wage pattern. The start and center rise. The end stays constant. The spread compresses.

These mechanics fit what nearshore engineering teams experience in practice. Automation stabilizes late stage steps, but that stability makes early and mid stages more valuable because their choices influence the entire probability of success. When teams in Latin America work in real time with the U.S., this structure becomes even more visible. Time zone alignment exposes every handoff. A missed step upstream is felt immediately. Downstream AI tools can correct some errors, but they also soften the discipline margin if used too aggressively. The model explains why hybrid approaches work, why humans remain central in architectural and integration roles, and why late stage automation improves reliability only when upstream incentives remain intact.

The path forward is clear. AI should handle the end of the chain where incentives are flat. Humans should anchor the middle where context and judgment matter most. The first position benefits from support but not full replacement. Organizations should expect wage compression as automation grows. They should also expect that leaving some human presence in the loop is not waste but a mechanism for preserving system wide discipline.

This framework does not attempt to predict every organizational nuance. It shows the structural forces that govern how people behave when automation reshapes the pipeline around them. It offers a map that

CTOs can use when designing nearshore delivery systems with heterogeneous cognitive roles, AI augmented workflows, and time aligned engineering teams. The model gives direction. Axiom Cortex and the Nearshore IT CoPilot supply the data. Together they make it possible to design distributed teams that remain reliable even as automation becomes a routine part of the engineering process.

This is the economic story of where AI fits, where humans still matter, and why the structure of incentives will continue to shape successful nearshore pipelines long after the tools themselves evolve.

## Glossary of Terms

### **AI unit**

A deterministic effort agent that always chooses effort equal to 1 and incurs fixed cost  $c$ . Represents automated reliability in the pipeline.

### **Axiom Cortex**

TeamStation AI's cognitive evaluation engine that generates neuro psychometric signals, reliability profiles, and behavioral gradients for engineering candidates and active team members.

### **Effort cost ( $c$ )**

The disutility to a human worker from choosing effort. Constant in the baseline model but may vary by role or cognitive load in advanced extensions.

### **Effort choice ( $e_i$ )**

The action of worker  $i$ . Effort equals 1. Shirk equals 0.

### **Expected project success ( $p_k$ )**

The probability the project succeeds when exactly  $k$  workers exert effort. Increasing in  $k$ . Satisfies strict complementarity.

### **Incentive margin ( $p_n$ minus $\zeta_i^x$ )**

The difference in success probability between worker  $i$  exerting effort and shirking. Drives the wage equation.

### **Nearshore IT CoPilot**

TeamStation's operational orchestration layer, which supplies the human telemetry, process feedback, and reliability markers referenced throughout this research.

### **Peer monitoring**

The informational structure where each worker observes only the effort of the previous worker. Captures the pipeline discipline that nearshore teams use implicitly.

### **Replacement policy ( $x$ )**

The vector  $x = (x_1, \dots, x_n)$  where  $x_i$  is the probability that worker  $i$  is replaced by an AI unit. Total AI capacity is constrained.

### **Sequential team**

A workflow in which each position hands a stage of work to the next. Used in the model because many engineering pipelines approximate this structure.

### **Wage $w_i^x$**

Success contingent payment promised to worker  $i$  under policy  $x$ . Derived from incentive compatibility.

### **$\zeta_i^x$**

Probability the project succeeds when worker  $i$  shirks under the AI placement policy  $x$ . Central to the model because all incentive distortions flow through it.

# **Mathematical Definitions and Notation**

### **Workers**

There are  $n$  human workers indexed by  $i = 1, \dots, n$ . Each chooses effort

$$e_i \in \{0, 1\}$$

with cost  $c$  when  $e_i = 1$ .

### **Effort profile**

A vector  $e = (e_1, \dots, e_n)$  that describes the team's action profile.

### **Success function**

$p_k$  = probability of project success when exactly  $k$  workers exert effort.

Increasing and strictly complementary:

$$p_{k+2} - p_{k+1} > p_{k+1} - p_k.$$

### **Replacement policy**

$x = (x_1, \dots, x_n)$  with

$$0 \leq x_i \leq 1$$

and

$$\text{sum of } x_i \leq 1.$$

AI replaces worker  $i$  with probability  $x_i$ .

### **Shirking success term**

For each worker  $i$  under policy  $x$ , define:

$$\zeta_i^x = \text{probability the project succeeds if worker } i \text{ chooses } e_i = 0.$$

This term depends on where AI is installed downstream, since AI increases the probability of rescue.

### **Incentive compatibility constraint**

Worker  $i$  exerts effort if:

$$p_n w_i - c \geq \zeta_i^x w_i$$

Solving for the minimum wage:

$$w_i^x = c \text{ divided by } (p_n - \zeta_i^x)$$

This is the core equation of the model.

### Expected cost of the principal

For each position  $i$ :

$$\begin{aligned} \text{expected cost contribution} &= \\ x_i c + (1 - x_i) p_n w_i^x & \end{aligned}$$

Total cost is the sum across  $i$ .

### Derivative with respect to $x_i$

$d/dx_i$  (Expected cost) decomposes into:

1. direct cost saving
2. direct incentive distortion
3. indirect incentive distortion propagated upstream

The optimal  $x_i$  sets this derivative equal to zero for interior solutions.

### Optimal replacement

Denoted  $x^*$ . Satisfies:

$\text{ExpectedCost}(x^*) \leq \text{ExpectedCost}(x)$  for all feasible  $x$ .

### Characterization of $x^*$

- high probability at  $i = n$
- moderate probability at  $i = 1$
- lowest probability at interior (middle) positions
- often interior ( $0 < x_i^* < 1$ )

### Wage pattern under $x^*$

- $w_n$  stays fixed

- $w_1$  rises
- $w_{\text{middle}}$  rises the most
- compression emerges naturally

# Data, Methodology, and Provenance Statement

This research uses a hybrid methodology combining:

## **1. Formal economic modeling**

Sequential team structures, moral hazard, and complementarity conditions grounded in the literature on multi agent incentives and organizational design.

## **2. Internal telemetry from TeamStation AI**

All empirical interpretations, cognitive alignment mappings, reliability curves, and incentive observations draw exclusively from internal platform data generated by:

- the Nearshore IT CoPilot
- Axiom Cortex neuro psychometric engine
- candidate readiness evaluations
- delivery performance audits
- outcome tracking from engineering teams across Latin America
- client retrospectives and pipeline diagnostics

No external proprietary datasets are used at any point.

## **3. Controlled real world observation**

The model's implications are compared with patterns observed in:

- senior engineering clusters
- architecture and integration roles
- QA and data verification stages

- multi timezone handoff structures
- AI augmented delivery pipelines
- candidate reliability distributions from Axiom Cortex

#### **4. Mathematical derivation and symbolic analysis**

All formulas, derivatives, comparative statics, and incentive structures presented in this paper come from analytic derivation. They are not approximations.

#### **5. Zero synthetic external data**

All qualitative examples reflect real operational patterns observed inside the TeamStation ecosystem.

All quantitative statements refer only to symbolic model parameters.

## **Legal and Intellectual Property Notice**

This paper and all underlying models, formulas, diagrams, incentive structures, cognitive mappings, and interpretive frameworks are protected by copyright.

Copyright © 2025 TeamStation AI. All rights reserved.

TeamStation AI, Axiom Cortex, Nearshore IT CoPilot, Nebula Talent Graph, and all related marks, system names, evaluation methods, and platform descriptions are proprietary to TeamStation AI. These names and systems are commercial assets, and unauthorized use in research, commercial settings, derivative works, or competing products is strictly prohibited. Any use of these terms in academic citations must reference this paper and the originating TeamStation AI research division.

All mathematical expressions, analytic derivations, and structural models presented here were developed internally by TeamStation AI. They are original research contributions. Redistribution, adaptation, training of machine learning models on this content, or inclusion in any automated system without written permission is not allowed.

All data interpretations, behavioral gradients, and incentive mappings in this paper rely exclusively on internal telemetry from the Nearshore IT CoPilot and Axiom Cortex. The findings are not based on external proprietary datasets. No client identifiable data is stored in this document. All internal signals represent aggregated, anonymized performance patterns from nearshore engineering teams across Latin America.

Nothing in this paper should be interpreted as legal, financial, or employment advice. The content is for research, education, and operational strategy within engineering organizations.

For licensing, permissions, or research collaboration, contact [lonnie@teamstation.io](mailto:lonnie@teamstation.io)