# [Transcript] AI Hustle: News on Open AI, ChatGPT, Midjourney, NVIDIA, Anthropic, Open Source LLMs / Bridging Human Insight with AI: A Conversation with Toloka AI CEO Olga Megorskaya

Welcome to the OpenAI podcast, the podcast that opens up the world of AI in a quick and concise manner.

Tune in daily to hear the latest news and breakthroughs in the rapidly evolving world of artificial intelligence.

If you've been following the podcast for a while, you'll know that over the last six months I've been working on a stealth AI startup.

Of the hundreds of projects I've covered, this is the one that I believe has the greatest potential, so today I'm excited to announce AIBOX.

AIBOX is a no-code AI app building platform paired with the App Store for AI that lets you monetize your AI tools.

The platform lets you build apps by linking together AI models like chatGPT, mid-journey, and 11Labs.

Eventually, we'll integrate with software like Gmail, Trello, and Salesforce so you can use AI to automate every function in your organization.

To get notified when we launch and be one of the first to build on the platform, you can join the waitlist at AIBOX.AI, the link is in the show notes.

We are currently raising a seed round of funding.

If you're an investor that is focused on disruptive tech, I'd love to tell you more about the platform.

You can reach out to me at jaden at AIBOX.AI, I'll leave that email in the show notes.

Welcome to the AI Chat podcast.

I'm your host, Jayden Shaper.

For today on the podcast, we have the incredible opportunity of being joined with Olga Magorska.

Olga is the founder and CEO of Toloka AI, which is a company at the forefront of driving AI development through human expertise.

With a background that includes a crucial role in data production at Yandex, she is an authority in fields like natural language processing, computer vision, but not just a business leader.

Olga is also an academic contributor.

She's co-authored research papers and speaking at renowned conferences like neural IPS and ICML.

Her work has also been featured in Forbes, VentureB, and a lot of other publications.

So thank you so much and welcome to the show today, Olga.

Hello.

Hello.

Nice to meet you.

Super happy to have you on the show.

The first thing I want to talk to you a little bit about and ask you is if you could give everyone a little, like a brief description essentially of what Toloka does and kind of what motivated you to start this company to begin with.

Yeah, sure.

So in general, we always say that AI and the modern technologies related to artificial

intelligence stay on three key pillars.

These are the models, the algorithms, the hardware, and the data.

And overall, our company is focusing on fixing everything related to the data for AI.

This relates to collecting training data to train AI models, collecting human insights in order to evaluate the quality of models and monitoring, moderating the outcomes of the models to make sure that they work in correspondence with ethical and compliance requirements of particular application.

Okay.

And I mean, take me back to like when you got this company started.

What was your inspiration?

Were you working somewhere seeing this big need?

What kind of got you started creating this company?

Yeah, actually, if we look maybe 10 years ago, that has always been a great irony of artificial intelligence, even though it is called artificial intelligence, it is quite often misviewed that the large part of AI production is actually relies on human insights.

For example, in order to train a self-driving car to drive across the streets, you need to feed the model behind the self-driving with literally millions of images on every image, every car, every pedestrian, every tree, every object is labeled and signed carefully by human.

Or if you train the search algorithm, you would need literally hundreds of thousands of pairs of users queries and the judgments of whether the found document was relevant or not in order to train the ranking algorithm to rank more relevant documents higher than least relevant.

Or if you are talking with a chatbot, how to evaluate whether the answers of the chatbot are relevant, spelled in the correct way, formulated in a human way, etc., etc., all these things actually require human efforts.

And that was the problem back like 10 years ago, even 5 years ago, even couple of years ago, that was one of the biggest bottlenecks in AI production is the part related to human labeling, because everything else is automated and hence easily scalable, but human operational part scales much harder.

You need specific skills and specific techniques, specific methodologies in order to manage human efforts in a scalable way.

And that is what we started with in Taloka, because for the long time we have been solving the problem of how to scale up human operations, human labeling to make it scalable on a really large amount of data.

And thus we ended up opening the crowdsourcing platform, where on the one hand every person can register as a performer, we call our performers Talokers, and on the other hand, every business or every manager, every AI scientist can register as a requester.

And the requester has posted their tasks for labeling on the platform, Talokers chose the tasks they are interested in and performed them.

And in between, we invested a lot into quality control and actually creating automated pipelines out of minor efforts of different people that would come up with reliable outcome of the

data sets.

That was the thing with which Taloka started, so we kind of invented the way how to use the efforts of millions of independent performers, millions of independent people in order to result into high quality and large scale training data set for AI.

That's amazing.

Very cool.

I'm curious, talking about using all these people and working on that, what kind of led you to emphasize the role of human insight in the development of AI and machine learning technologies is something you focus on?

So what led to that?

So in general, there are two parts.

One is purely, I would say, mechanical.

So in order to create the AI solution, one needed huge amounts of human labeled data by the way, this is the trend that is changing right now with the development of large language models.

We see that really substantial part of what you previously could have been done only with human efforts right now can be effectively done using large language models and we use it in our production as well.

So this part is right now being optimized much more and Taloka is now developing into the stage where we are actually pioneering in effective usage of large language models for labeling data.

But there is a second aspect which I think will not go anywhere and this is more of the ethical aspect of human judgment which becomes more and more important these days because we see right now with the raise of generative AI with large language models with such applications as ChatGPT with the rise of usage of API you can use ChatGPT in many, many different applications and services.

So that part of AI solutions generation becomes much more democratized but then there is the last mile which stays and in a way becomes much more important you ought to verify that actually these AI solutions work correctly that they are unbiased that they perform in the way you want them to perform.

And what we observe right now from looking at the industry and that's I think very interesting thing.

Previously very few companies really cared about the quality of their AI solutions.

If we are working with multiple different companies across the globe and I would say only the top most mature technological companies really cared about measuring the quality of their AI solutions.

But all the others took it like okay it works somehow and that's fine.

Right now we see that this approach is changing.

Right now much more the industry starts to take responsible AI seriously.

The industry starts to hear about the quality much more because in a way right now the models became such frighteningly convincing they can tell you something and you want to believe them.

So what the next thing I'm really curious about is how do you ensure the quality of data labeling when you're using these large language models.
This is really important.
We know that garbage in, garbage out, the quality of the content that you're putting in is really what makes a huge impact on getting something really quality out.
So how do you really ensure that those are high quality?
Yeah that's a very good question and well on the one hand there are lots of techniques and actually the majority of Taloka know how and patterns that we have in our products are related to quality control and in a way it's very interesting because some of the methods well overall it's a pure mathematics.
When you operate on the scale of large amounts of people who are making large amounts of effort in order to collect the ultimate data set you can apply a lot of statistical methods in order to predict the reliability and the expected accuracy of a given person at the level of his or her expertise.
You can mathematically count how many votes from which experts you need to aggregate in order to achieve the desired level of quality over the final label.
It's quite interesting that the first approaches to such mathematical approach towards data labeling were done back in 1979.
There is a model of David and skin that was invented in 1979, it's like 50 years.
In order to aggregate the votes and the judgments of different doctors in medical sphere when they needed to come to a conclusion about the stage of certain patient.
But then this model received the second birth in AI era and that's quite an interesting fact.
So on the one hand there are there is a lot of mathematics inside which you can apply in order to achieve high quality of the results.
On the second hand, and this is what we are observing more and more right now, like the development of AI in general means that algorithms can do more and more work than previously was done by human.
And this is also valid for data labeling as well.
And that means that humans and human judgments are needed in less amount of cases, but these cases require much higher level of expertise.
That's why where the whole industry is driving right now and where Telokia as a product is driving right now is towards expert labeling, which on the one hand combines all the previous techniques of quality control that we inherit from our previous stage of development of the product, but on the other hand it also requires working with specific dedicated experts.
And this the speed with which this trend is developing is actually amazing, like half a year ago we could not expect that right now we would be labeling data sets with the skills of senior software developers.
Right now, software developers are doing the job of data annotation in order to train AI algorithms for co-pilots, projects of different programming operators.
And that's the fascinating progress that is happening right now.
And the last part of it is actually, and that's quite interesting to observe.

I mentioned previously that overall industry is right now developing towards the understanding that quality of AI solution is something that matters.

And at the same time, it's the question of cultural DNA of a company, whether you care, whether you know and have this instinct to care about quality or not.

And in Telokia we definitely do have this in the DNA because we are trained to measure the quality of data annotation of data sets of AI solutions from the start.

And right now we apply the same methods and techniques in order to evaluate the quality of not only the data that labeled to train the models, but also to measure the upcoming quality of overall AI solution.

Very interesting.

This is obviously like a really useful piece of technology.

This is an incredible company here.

One question that I am curious about, and I'm sure a lot of people have, is when you're building and scaling this, what has been one of the biggest challenges you've had to overcome and how did you manage that?

Building this incredible company, it's so complex, there's so many technical aspects.

I'm sure there's been some a lot of difficult problems you've overcome, but talk to me about some of the biggest ones.

Well, probably every year you would ask me this question and every year I would answer different answers because we are facing, we're surviving from one challenge to another as probably any company and any startup.

I would say the last year, maybe a little bit less than one year, of course it is all about the new technological development and the rise of generative AI and large language models because on the one hand it requires from us certain technological transformation.

Just because I mentioned the more AI develops, the more job can be done by AI instead of humans and we know it better than anyone else and we are here running to be ahead of the curve in order to update our products in accordance to this technological development.

At the same time, it is super interesting to observe because we are kind of in the eye of the storm right now, everything that is going on in generative AI industry is passing by us because everybody who is developing certain solution needs data to training the algorithms to fine-tuning the generic models to their particular use case to do reinforcement learning with human feedback, the RLHF stage, they need also to monitor the quality to evaluate the metrics of the quality, etc.

So it is all passing in front of us and it is very interesting to observe the new trends, the new markets that are expected to raise with the rise of this new technological era.

So right now the biggest challenge is just to be adaptive as the time requires it.

Yeah, I think a lot of firms can relate with people on that, a lot of people are seeing all this new technology, I mean you have already been in it for a while but a lot of people are just seeing all these rapid advancements and trying to make sure they are staying on top of it and staying relevant and adapting to everything that is going on.

Something that I would love to ask you about is would you be able to talk a little bit about your own personal journey, what led you to work in technology, what led you to

Podtranscript.com

work in AI, like how did you come to a place where you are running this thing right now?

Yeah, my journey actually, let me recall how it started.

Overall, I studied mathematics in the university and then I've been working as an analyst, building statistical models, predicting different behaviors of, I actually worked in an investment bank, building some models of whatever probability or default, something like that.

Back then it was called statistics, then it started to be called machine learning, now it's called AI and etc.

So in a way, just as you mentioned before, on the one hand everything is changing so fast, on the other hand everything stays on the fundamentals that actually do not change that much.

So I started there but then at some moment I ended up in machine learning, actually with the need to fix this problem of lack of training data that was produced by humans and then I started the journey to actually solve this problem of scalability, of finding enough data to train and validate the machine learning models and apparently this task is the task that I am solving throughout all my professional life since then.

Tactics change because different technologies raise and fall but overall the general problem on the market stays the same because technology is developing but you need to provide tools and infrastructure to support this development.

Very cool.

So that brings us to today, you've told us a little bit about what you guys are building.

What I would be really curious about is what are some upcoming features or services that you're excited to roll out in the near future for Toloka?

Yeah, we're actually very excited to be rolling out our new platform which would allow businesses and engineers to actually deploy their models and then fine tune them with the help of the data labeling techniques that we have in order to apply them to their certain applications.

So how we foresee the area of our future development, like overall we foresee the world as the world where there are several really large foundational models that lie in the foundation of large variety of different GNI applications but then in order to apply the generic technology to the particular use case you would need some additional efforts to fine tune the generic model to apply it to your particular use case and then to measure the quality of the solution and constantly monitor that it works in accordance to your quality requirements and etc.

So we want to make sure that Toloka is a place where you can support all this life cycle of the GNI solutions with the help of all the layers of our infrastructure that we have.

That's incredible.

That's going to be very exciting.

I'm sure a lot of people are excited for that.

I'm sure you're excited to have that launch and I'm sure that's a huge project.

One question that I'm sure a lot of people are asking or thinking right now or that I would like to ask you is based off of everything you've learned in your career and working in your space, what's one piece of advice that you feel like you could give to people or companies that are looking at implementing AI to do different use cases within their company?

What are some areas they should focus on?

What's a piece of advice that you could give them?

Well, it probably would not be the piece of business advice but rather my wish towards the industry is to really care about the quality of what you are developing because AI is democratizing so it becomes easier and easier and easier to come up with some idea and to launch it into production. The biggest thing that I would hope to see in the development of the industry is being responsible about caring of what you're actually deploying to the world.

That is some amazing advice and really an incredible thing to wish for the entire industry. I think building things we're proud of is something that really needs to be focused on in AI so I really do appreciate that. If people want to get in contact with you or if people want to find out more about what you guys are building, where can they find that?

All right, well I will leave a link to your website in the show notes but thank you so much for coming on the podcast today. For the listeners, thanks so much for tuning in to the AI Chat podcast. Make sure to rate us wherever you listen to your podcasts if you want more phenomenal guests like this and have an amazing rest of your day. If you are looking for an innovative and creative community of people using ChatGPT, you need to join our ChatGPT creators community. I'll drop a link in the description to this podcast. We'd love to see you there where we share tips and tricks of what is working in ChatGPT. It's a lot easier than a podcast as you can see screenshots, you can share and comment on things that are currently working so if this sounds interesting to you, check out the link in the comment. We'd love to have you in the community. Thanks for joining me on the open AI podcast. It would mean the world to me if you would rate this podcast wherever you listen to your podcasts and I'll see you tomorrow.