



INTE
RVIE
W(H
TTP
S://R
ESE
ARC
H.RE
DHA
T.CO
M/B
LOG
/ART
ICLE
_CA
TEG
ORY
/INT
ERVI
EW/
)

Power surge: the push for sustainability

in high- performance computing and AI workloads

by Parul

Singh(<https://research.redhat.com/blog/article-author/parul-singh/>)

Han

Dong(<https://research.redhat.com/blog/article-author/han-dong/>)

Shaun

Strohmer(<https://research.redhat.com/blog/article-author/shaun-strohmer/>)

Red Hat Research Quarterly

November
2024(<https://research.redhat.com/blog/issue/november-2024/>)

INTERVIEW([HTTPS://RESEARCH.REDHAT.COM/BLOG/ARTICLE_CATEGORY/INTERVIEW/](https://research.redhat.com/blog/article_category/interview/))

Power surge: the push for sustainability in high-performance computing and AI workloads



(<https://research.redhat.com/blog/article/power-surge-the-push-for-sustainability-in-high-performance-computing-and-ai-workloads/>)

John Goodhue has perspective. He was there at the birth of the internet and the development of the BBN Butterfly supercomputer, and now he's a leader in one of the toughest challenges of the current age of technology—sustainable computing. Comparisons abound: one report says carbon emissions from cloud computing equal or exceed emissions from all [...]

by **Parul Singh**(<https://research.redhat.com/blog/article-author/parul-singh/>)

Han Dong(<https://research.redhat.com/blog/article-author/han-dong/>)

Shaun Strohmer(<https://research.redhat.com/blog/article-author/shaun-strohmer/>)

Search articles 🔍



11 min read

John Goodhue has perspective. He was there at the birth of the internet and the development of the BBN Butterfly supercomputer, and now he's a leader in one of the toughest challenges of the current age of technology—sustainable computing. Comparisons abound: one report says carbon emissions from cloud computing equal or exceed emissions from all commercial flights combined. Another suggested that by 2027 AI workloads could be using as much energy as a country the size of the Netherlands. That leaves computing and research communities with a conundrum: how do we solve real-world global challenges without worsening climate change or devastating power grids?

RHRQ asked Han Dong and Parul Singh, engineers working on open source projects related to energy efficiency, to talk with John about his work as director of the Massachusetts Green High-Performance Computing Center (MGHPCC), a joint venture of Boston University, Harvard, MIT, Northeastern, and the University of Massachusetts system. The MGHPCC

provides computing and storage resources for over 20,000 faculty and student researchers and educators, including the [MOC Alliance-supported New England Research Cloud \(NERC\)](https://research.redhat.com/blog/research_project/mass-open-cloud/) (https://research.redhat.com/blog/research_project/mass-open-cloud/) and the [Red Hat Collaboratory at Boston University](https://www.bu.edu/rhcollab/) (<https://www.bu.edu/rhcollab/>), and works to maximize energy efficiency while minimizing environmental impacts. They offer us a deep dive into the factors that bear on sustainability, from renewable energy sources and hardware choices to scheduling and tuning policies. The latter two, by the way, are the subject of Han and Parul's article on the PEAKS project, also in this issue.

—[Shaun Strohmer](https://research.redhat.com/blog/article-author/shaun-strohmer/) (<https://research.redhat.com/blog/article-author/shaun-strohmer/>), Ed.

Parul Singh: Tell us about your background in hardware. How did that lead you to eventually running the MGHPCC?

John Goodhue: Coming out of school, I ended up spending equal time on networks and computing. It was an exciting time for networking technology, with the technologies that became the foundations for the Internet just coming into play. At the same time, the company I was working at (BBN) was working on tiny networks that fit in a cabinet, interconnecting a large number of microprocessors to form a single computer system. That work led to the [BBN butterfly](https://en.wikipedia.org/wiki/BBN_Butterfly) (https://en.wikipedia.org/wiki/BBN_Butterfly), the second or third generation of the thing that we now call a cluster.

In 2010, the world was just beginning to

I ended up at the MGHPCC after leaving a startup in the high-performance computing business. It's obviously a very different thing:

*realize that energy
efficiency in
datacenters might
matter*

you can think of it as level minus one in the OSI (Open Systems Interconnection) stack. At the time we started to build it in 2010, the world was just beginning to realize that energy efficiency in datacenters might

matter. Before then, even the idea of turning your air conditioners off in the winter was kind of a “why bother.” That’s how far the dark ages people were. We drove pretty hard on energy efficiency, and our timing was kind of lucky that way.

Parul Singh: From your perspective, what got people to start thinking about sustainability in computing?

John Goodhue: There was a turning point in the thinking in the industry around 2005, exacerbated by the recession and also influenced by growing awareness of the need to pay attention to climate change. An early challenge was finding ways to work with our architects and engineers to take steps that seemed radical at the time but are fairly commonplace. Though even today, developers of datacenters are under great pressure to get the thing built and into operation so they can start to get return on investment. Being thoughtful about energy efficiency and environmental footprint often takes a back seat as a result.

Going forward, we are beginning to see people think about what’s going to happen when three

*Developers of
datacenters
are under
great pressure
to get the
thing built and
into operation
so they can
start to get
return on
investment.
Being
thoughtful*

things converge. First, you're seeing a dramatic increase in the amount of power consumed by datacenters, which is estimated to be as much as 3% of total worldwide energy demand, and may double in the coming years. Second, you're seeing large segments of the rest of the economy pivoting to electricity. The transportation sector and electric cars are a good example, but it's happening everywhere. Third, there's the introduction of renewables, which will make supply more intermittent.

about energy efficiency and environmental footprint often takes a back seat as a result.

The MGHPCC is lucky—maybe there's a combination of skill and luck—that we settled ourselves in a city (Holyoke, MA) whose municipal electric company delivers 100% green energy to us through hydroelectric power. But they're also forward-looking and interested in creative ways to make sure that the energy supply we receive has as light a footprint on the grid as possible. If you look out maybe 10 years, you may see, for example, battery storage as a factor. Our load is steady 24/7, but maybe we can lighten that load during peak hours and then make it heavy during non-peak hours in a managed way that potentially both saves money and offsets the impact of renewables. Tracking those types of change is a challenge that all datacenter operators face.



The Massachusetts Green High Performance Computing Center (MGHPCC) in Holyoke, MA

Han Dong: One thing we don't really know in the systems field is how power purchase agreements work with datacenters. Can you talk about that?

John Goodhue: The energy markets are a fascinating thing. They were developed when deregulation happened and transmission, generation, and delivery were split into three different things. The MGHPCC is dealing with the municipal electric company, which was grandfathered into being able to do all three. You don't get a discount for consuming more—or less, for that matter. They currently encourage us to have a load that's steady with the ability to drop or decrease if there's advantage to doing so.

If you look at the transmission rate, it's based on monthly peak consumption. There's an advantage to shaving peaks that accrues to our supplier. We need to have a piece of the agreement that incentivizes us to do that. If you look at something called the [forward capacity charge](https://www.iso-ne.com/markets-operations/settlements/understand-bill/item-descriptions/forward-charge)([https://www.iso-ne.com/markets-operations/settlements/understand-bill/item-descriptions/forward-](https://www.iso-ne.com/markets-operations/settlements/understand-bill/item-descriptions/forward-charge)

[capacity-market-fcm-daily-charge](#)), that's based on your annual peak usage. That leads to a double incentive during one month of the year, but nobody knows what month that is until the end of the year. If you go to the [ISO New England website\(https://www.iso-ne.com/\)](https://www.iso-ne.com/), there is endless material on how the markets work. It's remarkable that the grid and energy markets work as well as they do given the complexity of the market and the engineering that supports them.

Han Dong: You mentioned the MGHPCC is 100% renewable. What does that mean?

John Goodhue: I like to say there are three ways of viewing it. First, where do the electrons come from? A quantum physicist would say who knows—everywhere. The more practical view is the electrical engineering view. There is a dam on the Connecticut River that generates way more energy than what we consume. It is connected to the same power substation that our primary power feed is connected to. From an electrical engineering point of view, you can't get any more tightly coupled than that to a renewable source. And there's the added benefit that we're not relying on transmission lines that have other environmental impacts. It's all right there within a half mile. The third view is the market view, where renewable energy like hydro or solar energy carries a premium. So we also pay the price premium for renewable energy. At both an electrical engineering level as well as the market level, we are 100% carbon free.

Han Dong: Is it reasonable to consider renewables free energy?

John Goodhue: It is absolutely not free. There are capital costs and maintenance costs. The only difference is that you're not burning fuel from Texas or Pennsylvania or the Middle East. The upfront costs can be higher. The good news is that solar is getting to a point where the lifetime cost of



The Holyoke Dam as seen from South Hadley, MA, during the “freshet,” or spring thaw. Photo by Simtropolitan, CC BY-SA 3.0

delivering a kilowatt hour of renewable energy continues to decline, but “free” is not a description you would attach to any energy source.

Han Dong: When you say upfront costs, are you talking about the embodied carbon(<https://www.epa.gov/greenerproducts/what-embodied-carbon#:~:text=Embodied%20carbon%E2%80%94also%20known%20as,states%20of%20a%20product's%20life.>) that’s attached to some of these sources that we don’t really think about?

John Goodhue: There are things that emit carbon and there’s things that don’t, and the money the MGHPCC pays for electricity generation goes to places that don’t emit carbon.

Han Dong: At an NSF Workshop last month, one of the issues discussed was embodied carbon in terms of datacenter servers. The standard longevity, at least for hyperscalers, is that these servers last at most about

six years and then they replace them. There's a bigger push towards prolonging the life of some of this older hardware. How do we do that?

John Goodhue: That happens at a few different levels. First, lifetimes started to increase around 2005. Before then, CPU clock rates were steadily increasing from about the late 1980s through the early 2000s, and the performance improvements in every three-year cycle were significant enough to justify replacing the server.

That equation changed as clock rates hit the wall at about three megahertz. As a result, lifetimes are now closer to five and six years—at least that's what we've been seeing. GPUs are an exception, as architectural changes introduce compelling improvements every two years or so. But even there you can keep the older ones around, to support less performance-sensitive workloads for them, and the motherboard, at least, has to change.

Han Dong: What are some examples of the architectural changes that are driving GPU performance?

John Goodhue: Many of them have to do with AI workloads, which can use fixed-point arithmetic, and algorithm improvements have made it possible to use much lower predictions. It's a complete U-turn compared to the drive for better floating-point performance needed for the scientific simulation applications.

Parul Singh: When I was working at Boston University, I learned that the MGHPCC datacenter is cooled naturally. Could you say more about how that's done, using water and the cold temperatures in the New England area?

John Goodhue: I'll unpack that a little bit. The cooling system for the datacenter really operates in three stages. We circulate cold water—not extremely cold, about 65 degrees Fahrenheit—into the computer room and run it through heat exchangers that either remove heat from the air that's ejected by the servers or remove heat through water that circulates right next to the chips. The second stage transfers heat to a different water loop that circulates through a set of cooling towers. The cooling tower uses two kinds of processes to remove heat from the water that circulates through them. One is evaporation—evaporation is a cooling process, as we learned in elementary school—and the other is, say in January, it's just cold outside so the water cools down.

The chillers are the biggest energy hogs in the building, so we operate to minimize the use of that resource.

There are two ways of moving the heat from the computer center loop to the cooling tower loop. Most of the year, we just use a heat exchanger because the water from the cooling tower loop is cold enough. There are times during the year in Massachusetts when it's too hot outside to allow us to cool the

water enough using just heat exchangers. In that case, we use chillers, which are just refrigeration units that move heat from one place to another. The chillers are the biggest energy hogs in the building, so we operate to minimize the use of that resource. As I mentioned earlier, most datacenters didn't bother to minimize the use of chillers until the early 2000s, when datacenters started to get more energy conscious.

Han Dong: None of this information you're talking about is exposed to a user. Is there any value in getting this information to an end user? If there's a way to let someone know before they deploy their job that the cooling costs are likely to be high at a given time, is it worth incentivizing them to

think, “Maybe I shouldn’t run my heavy workload today because of the potential cooling costs?”

John Goodhue: I’ll start with an analogy. I’ve spent a lot of time building networks, and one of the key reasons the internet works at all is layering, right? When I send a packet with a source and a destination address in a checksum using IPV4, I don’t need to know anything about what happened next except whether the packet was dropped. Even that layer doesn’t tell me how to figure it out. There’s the layer above that, which works on end-to-end reliability when it’s needed, or maybe it’s tolerant of dropping packets on the way up the stack in the session that we’re running. If every layer had to know what was happening in every other layer, you could never change it and it would frequently break.

The analogy isn’t perfect but it holds. There are signals the grid can give our energy supplier that our energy supplier can pass through to us about peak demand, for example, and backing off when a multi-peak happens. It’s feasible to pass those signals through, but you would need to react quickly. There is some promising research at BU and other universities looking at how to do that while minimizing impact on performance.

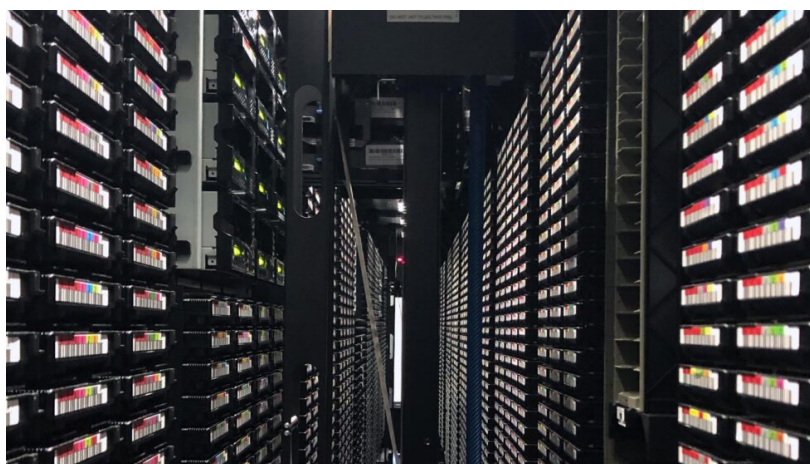
It is also possible to react to a signal at a lower level in the stack, by lowering the CPU clock rate. This has the advantage of being invisible to the application software, other than a decrease in performance.

Han Dong: The operating system has mechanisms to do exactly what you’re talking about. It just doesn’t have the policies there right now.

Let’s address the question everyone’s concerned about. From your perspective as director of a high-performance computing facility, what are the biggest challenges presented by generative AI and LLMs? For

example, in terms of workload, how much of the work in MGHPCC consists of things like generative AI and large language models? I'm assuming it's increasing, right?

John Goodhue: It is increasing. If I use GPU installations as a proxy, it's growing substantially. The reason I can't say exactly is there's been a wave of installations with people who put in orders last year but supply chain problems kept them from actually getting their stuff until the last three months or so, so the systems are just beginning to come up.



Racks of the Northeast Storage Exchange, created to address the escalating need for research data storage, housed at the MGHPCC

Han Dong: Are a lot of them doing training jobs right now or are they also setting up inferencing servers on the MGHPCC?

John Goodhue: That's a good question. I suspect that there tends to be more training, compared to industrial workflows, but I don't know what the mix is.

Han Dong: Are you thinking about how the MGHPCC might adapt? Do you have to change how you think about the cooling and layout of the datacenter?

John Goodhue: The MGHPCC has been able to take this change in stride, in part due to early design decisions, and in part because we started supporting systems with high power density in 2016, well before the AI wave appeared.

Not much has changed with the AI boom, with one exception, which is the amount of power that gets consumed per square foot. Ten years ago, high-end enterprise workloads were six kilowatts per rack (KPR), and research computing workloads were probably more like 12 KPR. GPU-heavy research computing workloads can now be as much as 60 or 70 KPR, going up to 100 KPR. So we are seeing the amount of computing resources per square foot increase by factors of two and three. How we distribute power hasn't changed our cooling much at all.

More broadly, the [Coalition for Academic Scientific Computation \(CASC\)](https://casc.org/researchpub/position-statements/) (<https://casc.org/researchpub/position-statements/>) has organized several working groups around the topics of AI, energy efficiency, and building and maintaining datacenters like the MGHPCC. The technology is evolving rapidly, and it's essential that research computing and datacenters do not fall behind their for-profit counterparts.

Parul Singh: Thanks for your time, John. This has been quite interesting and very helpful.

The technology is evolving rapidly, and it's essential that research computing and datacenters do not fall behind their for-profit counterparts.

SHARE THIS ARTICLE

MORE LIKE THIS

INTERVIEW

From silos to startups: why universities must be part of industry's AI growth(<https://research.redhat.com/blog/article/from-silos-to-startups-why-universities-must-be-part-of-industrys-ai-growth/>)

Brian Stevens

INTERVIEW

A marriage of true minds: Making university-industry collaborations succeed(<https://research.redhat.com/blog/article/a-marriage-of-true-minds-making-university-industry-collaborations-succeed/>)

Martin Ukrop

INTERVIEW

AI DIY: How research is making custom language models work with more of us(<https://research.redhat.com/blog/article/ai-diy-how->

research-is-making-custom-language-models-work-with-more-of-us/)

Heidi Dempsey

“How many lives am I impacting?” That’s the question that set Akash Srivastava, Founding Manager of the Red Hat AI Innovation Team, on a path to developing the end-to-end open source LLM customization project known as InstructLab. A principal investigator (PI) at the MIT-IBM Watson AI Lab since 2019, Akash has a long professional history [...]

INTERVIEW

ChRIS five years later: the groundbreaking platform levels the playing field for advanced analytics and AI in medicine(<https://research.redhat.com/blog/article/chris-five-years-later-the-groundbreaking-platform-levels-the-playing-field-for-advanced-analytics-and-ai-in-medicine/>)

Orran Krieger

Shaun Strohmer

What if there were an open source web-based computing platform that not only accelerates the time it takes to share and analyze life-saving radiological data, but also allows for collaborative and novel research on this data, all hosted on a public cloud to democratize access? In 2018, Red Hat and Boston Children’s Hospital announced a [...]

INTERVIEW

Future vision: on the internet, technopanic, and the limits of AI(<https://research.redhat.com/blog/article/future-vision-on-the-internet-technopanic-and-the-limits-of-ai/>)

Jason Schlessman

Everyone has an opinion on misinformation and AI these days, but few are as qualified to share it as computer vision expert and technology ethicist Walter Scheirer. Scheirer is the Dennis O. Doughty Collegiate Associate Professor of Computer Science and Engineering at the University of Notre Dame and a faculty affiliate of Notre Dame's Technology [...]

INTERVIEW

No more gatekeepers: Why technological ignorance is radically dangerous and how an open world will help(<https://research.redhat.com/blog/article/no-more-gatekeepers-why-technological-ignorance-is-radically-dangerous-and-how-an-open-world-will-help/>)

Jason Schlessman

What is the role of the technologist when building the future? According to Boston University professor Jonathan Appavoo, "We must enable flight, not create bonds!" Professor Appavoo began his career as a Research Staff Member at IBM before returning to

academia and winning the National Science Foundation's CAREER award, the most prestigious NSF award for [...]

INTERVIEW

"Research is an adventure": Putting theory to the test at the university and in the field(<https://research.redhat.com/blog/article/research-is-an-adventure-putting-theory-to-the-test-at-the-university-and-in-the-field/>)

Martin Ukrop

Don't tell engineering professor Miroslav Bureš that software testing can't be exciting. As the System Testing IntelLigent Lab (STILL) lead at Czech Technical University in Prague (CTU), Bureš's work bridges the gap between abstract mathematics and mission-critical healthcare and defense systems. His research focuses on system testing and test automation methods to give people new [...]

INTERVIEW

"That's what open source is all about": A short history of collaboration, innovation, and education in research(<https://research.redhat.com/blog/article/thats-what-open-source-is-all-about-a-short-history-of-collaboration-innovation-and-education-in-research/>)

Shaun Strohmer

In 2017, Red Hat Chairman Paul Cormier and Boston University (BU) professor Orran Krieger helped spearhead a collaborative

partnership between the two institutions that would come to include expanding Red Hat's participation in the MOC Alliance, the establishment of the Red Hat Collaboratory at BU for research incubation, and the creation of a Red Hat [...]

INTERVIEW

Where are we with wireless? How researchers are pushing forward the state of the art, and what that means for industry(<https://research.redhat.com/blog/article/where-are-we-with-wireless-how-researchers-are-pushing-forward-the-state-of-the-art-and-what-that-means-for-industry/>)

Heidi Dempsey

ABOUT THE AUTHOR



Parul Singh

Parul Singh is a software engineer in the Office of the CTO at Red Hat. She is currently leading Red Hat's effort on scaling ChRIS using OpenShift. Her research lies in AI, ML, NLP, and big data.

ABOUT THE AUTHOR

**Han Dong**

Han Dong is a postdoc in the Computer Science department at Boston University. His research interests lie in distributed systems, high-performance computing, and operating systems. He is interested in research addressing the growing energy needs of our modern systems.

ABOUT THE AUTHOR

**Shaun Strohmer**

Shaun Strohmer is the editor of the *Red Hat Research Quarterly*. She has worked as a writer and editor in academic publishing for over twenty years, and since 2014 she has focused on software development, cybersecurity, and computer science.

John Goodhue

RELATED PROJECTS

- [Mass Open Cloud \(MOC\): An open, distributed platform enabling AI/ML workloads](https://research.redhat.com/blog/research_project/mass-open-cloud/)(https://research.redhat.com/blog/research_project/mass-open-cloud/)

ARTICLE FEATURED IN

(<https://research.redhat.com/blog/issue/november-2024/>)
(<https://research.redhat.com/blog/issue/november-2024/>)

Red Hat Research Quarterly

November 2024

Download PDF 

Subscribe now

IN THIS ISSUE

FROM THE DIRECTOR

From particles to prototypes: what we learn from managing open clouds(<https://research.redhat.com/blog/article/from-particles-to-prototypes-what-we-learn-from-managing-open-clouds/>)



Heidi Dempsey

NEWS

**Observability cluster added to the MOC Alliance's
New England Research**

Cloud(<https://research.redhat.com/blog/article/observability-cluster-added-to-the-moc-alliances-new-england-research-cloud/>)



Thorsten Schwesig



Christopher Tate

NEWS

Publication highlights—November

2024(<https://research.redhat.com/blog/article/publication-highlights-november-2024/>)

INTERVIEW

Power surge: the push for sustainability in high-performance computing and AI workloads(<https://research.redhat.com/blog/article/power-surge-the-push-for-sustainability-in-high-performance-computing-and-ai-workloads/>)



Parul Singh



Han Dong



Shaun Strohmer

FEATURE

Scaling the PEAKS of sustainability with insights from Kepler and machine learning(<https://research.redhat.com/blog/article>

/scaling-the-peaks-of-sustainability-with-insights-from-kepler-and-machine-learning/)



Han Dong



Parul Singh

FEATURE

The Open Education Project is ready to scale(<https://research.redhat.com/blog/article/the-open-education-project-is-ready-to-scale/>)



Danni Shi

FEATURE

Open source authentication exposed: how open source developers perceive user authentication(<https://research.redhat.com/blog/article/open-source-authentication-exposed->

how-open-source-developers-perceive-user-authentication/)



Agáta Kružíková

COLUMN

Red Hat and the MOC-A: creating the open source cloud for the AI era(<https://research.redhat.com/blog/article/red-hat-and-the-moc-a-creating-the-open-source-cloud-for-the-ai-era/>)



Orran Krieger



Heidi Dempsey

LEARN

Research

Areas(<https://research.redhat.com/research/>)

Masters'

Theses(<https://research.redhat.com/research/theses/>)

Events(<https://research.redhat.com/events/>)

News(<https://research.redhat.com/news-2/>)

Magazine(<https://research.redhat.com/quarterly/>)

ENGAGE(<https://research.redhat.com/get-involved/>)

Contact

Us(<https://research.redhat.com/feedback/>)

Log In(<https://research.redhat.com/wp-login.php>)

ABOUT

Red Hat Research connects Red Hat engineers with professors, researchers, and students to bring great research ideas into open source communities. Our activities around the world have produced grants from government and industry, papers at top conferences, and results that have landed in open source projects of all kinds. Red Hat

Research welcomes participation from research-minded individuals around the world.

[f](#) [X](#) [@](#) [v](#) [in](#)

(htt ps:/ /ww w.fa
(htt ps:/ /twi tter.
(htt ps:/ /ww w.in
(htt ps:/ /ww w.yo
(htt ps:/ /ww w.lin
cebu
ook.
com
stag
utu
kedi
/Re
ram.
be.c
n.co

Copyright © 2025 Red Hat, Inc.
com dHa com om/ m/c
/Re t) /red user om
dHa hatj. /Re pan
Terms of
tInc nc/ dHa y/re
use(https://www.redhat.com/en/about/t
) ? tVid d-
erms-use)
hl=e eos) hat/
n))



Privacy
statement(https://www.redhat.com/en/a
bout/privacy-policy)
All policies and
guidelines(https://www.redhat.com/en/abo
ut/all-policies-guidelines)