

Modularer Rechenzentrumsbau

und was das mit GPU Servern und dem RZReg zu tun hat

Sebastian Krey



About GWDG



NHR-NORD@GÖTTINGEN



- IT service center and data center operation for **University Göttingen** and **Max Planck Society** (MPG) since 1970
- Operating site of “North German Supercomputing Alliance” (**HLRN**) since 2018, since 2021 part of **NHR**
- AI Service Center **KISSKI** for critical infrastructure
- HPC operating site for the “German Aerospace Center” (**DLR**) since 2022

Network for National High-Performance Computing

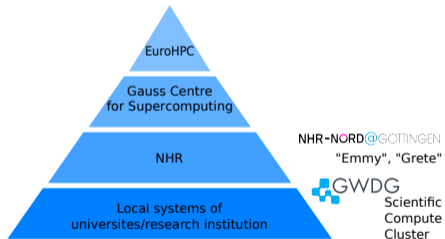
NHR-NORD@GÖTTINGEN

UGOE + GWDG

- Since 2021: Funding for national Tier-2 supercomputing (62.5M=C p.a.)
 - ▶ Nine centres
 - ▶ Annual funding 7,3M€ p.a.
- Usable for researchers at all German universities



HPC systems at GWDG

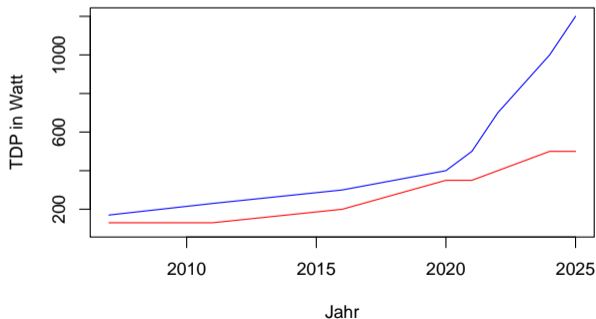


- Tier 2: **HLRN/NHR "Emmy"**
Top500 #47 Nov. 2020, now #219
- Tier 2: **NHR/KISSKI "Grete"**
Top500 #141 Nov. 2023, Green500 #22,
now #226/#70
- Tier 2: **NHR/KISSKI "Grete Phase 3"**
Top500 #274 Nov. 2024, Green500 #24,
now #315/#35
- Tier 3: **Scientific Compute Cluster**
- **"CARO" for DLR**
Top500 #135 Nov. 2021, now #335
- Several smaller systems for MPG and UGOE

Leistungsdichte durch KI

Zunehmende Benutzung von Beschleunigerarchitekturen zur Ausführung von KI Modellen in allen RZ Bereichen.

TDP Entwicklung von CPUs und GPUs seit 2007



HPC-sierung von RZ Workloads
→ über 100kW pro Schrank
“normal”.
→ Zunehmender Bedarf
(direkter) Wasserkühlung

EnEFG Anforderungen an RZ

- Verbesserung der Rechenzentrumseffizienz
 - ▶ Bestand: PUE <1.5 ab Juli 2027 und <1.3 ab Juli 2030
 - ▶ Neubauten ab Juli 2026: PUE <1.2 (EER >5)
- Nutzung erneuerbarer Energien zur Stromversorgung (bilanziell 50% seit 2024, 100% ab 2027)
- Abwärmenutzung für Neubauten mit Inbetriebnahme ab
 - ▶ Juli 2026 min. 10%
 - ▶ Juli 2027 min. 15%
 - ▶ Juli 2028 min. 20%
- Energie- und Umweltmanagementsystem gemäß ISO 50001 (mit Zertifizierung ab 2026 für >1MW, Ausnahme <7.5GWh/a und >50% Abwärmenutzung)
- Informationspflichten zur Energienutzung gegenüber RZReg und Nutzer

Modernisierung vs Neubau

- Effizienzoptimierungen erfordern Modernisierung der RZ Infrastruktur
- Höhere Leistungsdichten erfordern Modernisierung der RZ Infrastruktur
- Bauen im Bestand kompliziert (Restriktionen des Gebäudes)
- Wie Uptime der bestehenden IT sicherstellen
- Neubau braucht viel Zeit, Platz, kein Bestandsschutz bei Effizienz und Abwärmenutzung
- HPC-sierung der klassischen IT macht langfristige Infrastrukturplanungen immer schwieriger
- Klassischer RZ Bau (Gebäude für >20 Jahre) noch zeitgemäß?
- Modularisierung für optimale Auslastung (>80%) der Infrastruktur?

RZReg Meldepflichten

Alle Rechenzentren ab 300 kW Anschlussleistung müssen Daten liefern.

- Einreichungsfrist 31. März für Daten des Vorjahres
- Inbetriebnahmedatum
- Größe des RZ und des Rechnerraums
- Anschlussleistung
- Adresse und Kontaktdaten
- Redundanzlevel aller Infrastrukturkomponenten (Gebäude- bis Rackebene)
- Eingesetzte Kältemittel
- Rackspace getrennt für Storage und Server
- Stromverbräuche
 - ▶ Gesamtes Rechenzentrum
 - ▶ IT
 - ▶ Kühlsysteme
 - ▶ Aufbereitung zur Abwärmenutzung

RZReg Meldepflichten

- Menge und Art erneuerbare Energien (selbst erzeugt, Power-Purchase-Agreement, Herkunftsnachweise)
- Abwärmenutzung
- Wasserverbräuche für IT Betrieb (gesamt und Trinkwasser)
- Storagekapazität
- Rechenkapazität (SPEC SERT) (gesamt und neu installierte Hardware)
- Netzwerkbandbreite
- Gesamttraffic
- Durchschnittliche Temperaturen von
 - ▶ Zuluft, Abluft
 - ▶ Wasser Vorlauf und Rücklauf (Temperaturniveau der Abwärme)

Monitoringdetails

Auskunftsanspruch einzelner Kunden des RZs zur Ressourcennutzung ihrer IT erfordert deutlich mehr Details als RZReg Meldung

- Ressourcenverbrauch der Gebäudeinfrastruktur getrennt
 - ▶ Kühlkreisen
 - ▶ Stromversorgungen
- Wärmemengen
- Leistungsaufnahmen aller Geräte (einzelne Pumpen, Server, Switches)
- Wassermengen (DLC, wassergek. Schränke/Rücktüren, Adiabatik)
- Datenmengen (Storage und Netzwerk)
- Energieeffizienz sowie Rechen- und Speicherkapazität neuer Hardware

Monitoring

- Verschiedenste Feldbusse und Protokolle
- Unterschiedliche Anforderungen an zeitliche Auflösung (Reporting und Betriebsoptimierung)
- Große Anzahl an Datenpunkten
- Zeitreihendatenbank z.B. Prometheus oder InfluxDB
- Visualisierung und Datenaggregation mit Grafana
- Datensammlung mit verschiedenen Exportern (Modbus, SNMP, IPMI, Redfish, etc.)
- Erfassung aller Geräte und Verkabelung in Nautobot als DCIM Tool

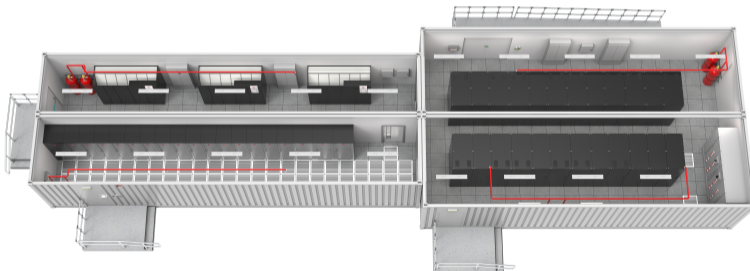
Standortsituation Uni Göttingen 2018

- 2017 Zuschlag für erstes Tier 2 HPC System (HLRN IV) am Standort
- Erstes Warmwasser gekühltes System am Standort
- Erste Installationsphase H2 2018, Vollausbau Ende 2019
- 2018 nach Änderung des Änderungsvertrags drastische Anforderungsänderungen für Rechenzentrumsinfrastruktur:
 - ▶ 1,32 MW anstatt 800 kW
 - ▶ 26-28 anstatt 12-13 Schränke
- Rechenzentrumsstandort der ersten Phase nicht ausreichend, Kostenvoranschläge für Umbau >5M EUR.
- Rechenzentrumsneubau in Arbeit aber nicht vor 2021 bereit

Lösungskonzept

- Bau eines Modulares Rechenzentrums
- Gelände vorhanden
- Sehr knappes Budget (<3M EUR)
- Eng getakteter Zeitplan, Zuschlag für MDC Bau im Juli 2019
- Enge Abstimmung mit Systemanbieter Atos bzgl. Infrastrukturanforderungen
- Zwei Räume:
 - ▶ 70qm für luftgekühlte Systeme (Phase 1, Storage, Managementserver): 300kW, 70kW USV, 19 Schränke, 10-26kW pro Schrank
 - ▶ 80qm für DLC Systeme: 1,1MW, 20% Restwärme über Seitenkühler, 96kW pro Schrank, CoolIT DLC

Gebäudeschnitt



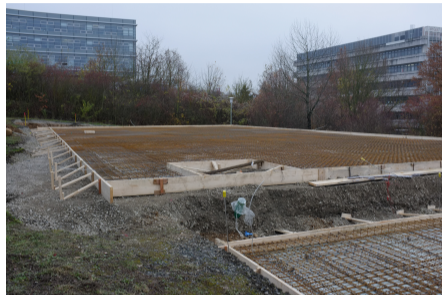
Kühlungsauslegung

- HPL Leistungsaufnahme außerhalb der Abnahme unrealistisch
- Alltagsleistungsaufnahme eher 70-80% von HPL
- Temperaturen in Göttingen moderat:
33°C sind 99,9% Quantil, 28,4°C sind 99%
- MDC Standort kein Wasseranschluss → trockene Rückkühler
- Weitere Anforderungsänderung im Herbst 2019:
32°C Vorlauf anstatt 36°C
- Zulufttemperaturen so hoch wie effizient möglich

Ergebnis Kühlungsauslegung

- HPL Leistungsaufnahme bis 20°C Außentemperatur
- Bei 32°C müssen 75% HPL möglich sein
- Trockene Rückkühler mit $dT=4K$ zur Außentemperatur
- Gemeinsamer Kreislauf für Umluftklimageräte, Seitenkühler und DLC
- Außentemperatur geführte Vorlauftemperatur
- Hybride Umluftklimageräte und Seitenkühler

Fotos Aufbauphase Nov/Dez 2019



Fotos Aufbauphase 02.-04.06.2020



Fotos Aufbauphase 17.06.-09.07.2020



Fotos Aufbauphase 30.07.2020



Aufbauphase

- November 2019, Gießen der Bodenplatte, LWL und Bau Trafostation
- 3 Monate Verzögerung beim Bau der Containermodule
- Weitere 3 Monate durch pandemiebedingte Schließung der Grenzen
- Ankunft der Containermodule und Rückkühler ab 02.06.2020
- Juni 2020 Anbindung an Infrastruktur, Bau des Primärkreislaufs, Funktionstests
- 01.07.2020 Atos beginnt Umzug der ersten Installationsphase
- 25.08.2020 Phase 1 im MDC in Produktion
- 07.10.2020 DLC gekühlte Phase 2 bereit für 48h Lasttest

PUE Werte und Optimierungen

- Erstes Betriebsjahr PUE 1,13 gesamt und 1,07 für DLC Raum
- Feintuning der Zulufttemperatur auf 24,5°C für luftgekühlte System und 27,5°C für DLC Systeme
- DLC Temperaturen abhängig vom Gebäudevorlauf für höhere Heat Capture Rate bei kühlen Temperaturen
- Alle PID Regler zu langsam für große Lastsprünge, also nachjustieren
- Außentemperaturabhängige Drehzahlbegrenzung der Rückkühler
- 12-Monats PUE seitdem 1,11-1,15, DLC Raum 1,05-1,08
- Schlechtester Monats PUE im August 2022 mit 1,20 bzw. 1,10

Zusammenfassung

- Leistungsdichte der Prozessoren steigt massiv
- (Direkte-) Wasserkühlung wird immer wichtiger
- Detailliertes Monitoring für Berichtspflichten notwendig
- Beides erfordert massive Modernisierung der RZ Infrastruktur
- Übergeordnetes Monitoring mit Open Source Produkten
- Modulare Rechenzentrumsbauweise hat sich bewährt (Geschwindigkeit und Kosten)
- Trotz technisch einfachen Aufbaus sehr gute Energieeffizienz
- Einfacher Aufbau erleichtert Feintuning der Komponenten
- Bis 40°C Außentemperatur Betrieb mit trockenen Rückkühlern möglich
- Zukunft: Abwärmenutzung in den Nachbargebäuden Physik und Chemie